

A THEORY OF NONPARAMETRIC REGRESSION IN THE PRESENCE OF COMPLEX NUISANCE COMPONENTS

MARTIN WAHL

ABSTRACT. In this paper, we consider the nonparametric random regression model $Y = f_1(X_1) + f_2(X_2) + \epsilon$ and address the problem of estimating the function f_1 . The term $f_2(X_2)$ is regarded as a nuisance term which can be considerably more complex than $f_1(X_1)$. Under minimal assumptions, we prove several nonasymptotic $L^2(\mathbb{P}^X)$ -risk bounds for our estimators of f_1 . Our approach is geometric and based on considerations in Hilbert spaces. It shows that the performance of our estimators is closely related to geometric quantities, such as minimal angles and Hilbert-Schmidt norms. Our results establish new conditions under which the estimators of f_1 have up to first order the same sharp upper bound as the corresponding estimators of f_1 in the model $Y = f_1(X_1) + \epsilon$. As an example we apply the results to an additive model in which the number of components is very large or in which the nuisance components are considerably less smooth than f_1 . In particular, the results apply to an asymptotic scenario in which the number of components is allowed to increase with the sample size.

1. INTRODUCTION

In this paper, we consider the nonparametric random regression model

$$Y = f_1(X_1) + f_2(X_2) + \epsilon. \quad (1.1)$$

We study the problem of estimating the function f_1 , while the function f_2 is regarded as a nuisance parameter. We are interested in settings where the second term $f_2(X_2)$ is much more complex than the first term $f_1(X_1)$. A particular model of interest is the additive model

$$Y = f_1(X_1) + \sum_{j=1}^{q-1} f_{2j}(X_{2j}) + \epsilon \quad (1.2)$$

2010 *Mathematics Subject Classification.* 62G08, 62G20, 62H05, 62H20.

Key words and phrases. Nonparametric regression, nuisance components, projection on sumspaces, minimax estimation, additive model, increasing number of components.

in which the nuisance components f_{2j} are considerably less smooth than f_1 or in which the number of components q is very large, for instance in the sense that q is allowed to increase with the sample size n . The estimation problem is similar to the one arising in semiparametric models where the aim is to estimate a finite-dimensional parameter in the presence of a (more complex) infinite-dimensional parameter.

Estimation in nonparametric additive models is a well-studied topic, especially when considering the problem of estimating all components in the case that q is fixed. One of the seminal theoretical papers is by Stone [30], who showed that each component can be estimated with the rate of convergence corresponding to the situation in which the other components are known. Since then, many estimation procedures have been proposed, many of them consisting of several steps. In the work by Linton [19] and Fan, Härdle, and Mammen [11], it is shown that there exist estimators of single components which have the same asymptotic bias and variance as the corresponding oracle estimators for which the other components are known.

Probably the most popular estimation procedures are the backfitting procedures, which are empirical versions of the orthogonal projection onto the subspace of additive functions in a Hilbert space setting (see, e.g., the book by Hastie and Tibshirani [12] and the references therein). This orthogonal projection was studied, e.g., by Breiman and Friedman [6] (see also the book by Bickel, Klaassen, Ritov, and Wellner [3, Appendix A.4]). They showed that, under certain conditions including compactness of certain conditional expectation operators, it can be computed by an iterative procedure using only bivariate conditional expectation operators. Replacing these conditional expectation operators by empirical versions leads to the backfitting procedures. Opsomer and Ruppert [24] and Opsomer [23] computed the asymptotic bias and variance of estimators based on the backfitting procedure in the case where the conditional expectation operators are estimated using local polynomial regression. Mammen, Linton, and Nielsen [20] introduced the smooth backfitting procedure and showed that their estimators of single components achieve the same asymptotic bias and variance as oracle estimators for which the other components are known. Concerning the distribution of the covariates, they make some high-level assumptions which are satisfied under some boundedness conditions on the one- and two-dimensional densities. This is still more than is required in the Hilbert space setting (see [6]). In the work by Horowitz, Klemelä, and Mammen [13], a general two-step procedure was proposed in which a preliminary undersmoothed estimator is based on the smooth backfitting procedure of [20]. They also showed that there are

estimators which are asymptotically efficient (i.e., achieve the asymptotic minimax risk) with the same constant as in the case with only one component. In addition to the assumptions coming from the results in [20], they require a Lipschitz condition for all components.

The problem of estimating f_1 in cases in which $f_2(X_2)$ is more complex than $f_1(X_1)$ is also considered in the work by Efromovich [10] and Muro and van de Geer [22]. In [10], an estimator of f_1 is constructed which is both adaptive to the unknown smoothness and asymptotically efficient with the same constant as in the case with only one component. The assumptions include smoothness and boundedness conditions on the full-dimensional density of (X_1, X_2) . The construction of the estimator is involved and starts with a blockwise-shrinkage oracle estimator. In [22], a penalized least squares estimator is analyzed in cases where the function f_1 is smoother than the function f_2 . Under certain assumptions including smoothness conditions on the design densities, it is shown that for both components, the estimator attains the rate of convergence corresponding to the situation in which the other component is known; i.e., no undersmoothing of the function f_2 is needed to estimate the function f_1 .

The previously discussed literature on additive models focuses on the asymptotic behavior of estimators as the number of observations n goes to infinity in the case that q is fixed. Note that one of our purposes is to generalize several results to the case that q increases with n .

Recently, high-dimensional sparse additive models have been studied, e.g., in the work by Meier, van de Geer, and Bühlmann [21], Huang, Horowitz, and Wei [14], Koltchinskii and Yuan [16], Raskutti, Wainwright, and Yu [25], Suzuki and Sugiyama [31], and Dalalyan, Ingster, and Tsybakov [7]. These papers consider the case that the number of covariates q is much larger than the sample size n . The focus is on the problem of estimating all components under sparsity constraints. In [7], e.g., the authors construct an estimator achieving optimal minimax rates of convergence. These rates of convergence depend on q and also on the smallest degree of smoothness of the f_{2j} . Hence, they may only lead to crude bounds for the rates of convergence of estimators of f_1 . Let us mention that in this paper, we do not consider a sparsity scenario. We are interested in cases in which the number of components q is very large, but smaller than n .

In this paper, we consider model (1.1) in the case that the functions f_1 and f_2 belong to closed subspaces H_1 and H_2 of $\{g_1 \in L^2(\mathbb{P}^{X_1}) : \mathbb{E}[g_1(X_1)] = 0\}$ and $L^2(\mathbb{P}^{X_2})$, respectively. We propose an estimator of f_1 which is based on the composition of two least squares criteria. Our main contribution is to derive several nonasymptotic risk bounds

which show that the performance of our estimators is closely related to geometric quantities of H_1 and H_2 , such as minimal angles and Hilbert-Schmidt norms. These risk bounds lead to minimal conditions under which the function f_1 can be estimated (up to first order) just as well as in the model $Y = f_1(X_1) + \epsilon$. Our analysis is based on geometric considerations in Hilbert spaces, and relies on the theory of projections on sumspaces in Hilbert spaces (see, e.g., [3, Appendix A.4]). Moreover, we apply recent concentration inequalities for structured random matrices (see, e.g., the work by Rauhut [26]) in order to show that several geometric properties in the Hilbert space setting carry over to the finite sample setting with high probability. As a main example we apply our results to the additive model (1.2) which corresponds to the case that H_2 has an additive structure. Using our results, we establish new conditions on q and on the smoothness of the nuisance components under which our estimator of f_1 attains the same (nonasymptotic) optimal rate of convergence as the corresponding least squares estimator in the model $Y = f_1(X_1) + \epsilon$. We also address the question of when the corresponding constants coincide.

The paper is organized as follows. In Section 2 and 3, we present the assumptions on the model and state our main results in Theorems 1-4. In Section 4, we apply our results to several models including the additive model. The proofs of our results are given in Sections 5 and 6. Finally, some complements are given in the Appendix.

2. THE FRAMEWORK

2.1. The model. Let (Y, X_1, X_2) be a triple of random variables satisfying (1.1), where X_1 and X_2 take values in some measurable spaces (S_1, \mathcal{B}_1) and (S_2, \mathcal{B}_2) , respectively, ϵ is a real valued random variable such that $\mathbb{E}[\epsilon|X] = 0$ and $\mathbb{E}[\epsilon^2|X] = \sigma^2$, and the unknown regression functions satisfy the following assumption:

Assumption 1. *Suppose that $f_1 \in H_1$, where*

$$H_1 \subseteq \{g_1 \in L^2(\mathbb{P}^{X_1}) : \mathbb{E}[g_1(X_1)] = 0\}$$

is a closed subspace, and that $f_2 \in H_2$, where $H_2 \subseteq L^2(\mathbb{P}^{X_2})$ is a closed subspace.

Structural assumptions on f_1 and f_2 (see, e.g., Section 4 where we also consider the additive model) should be incorporated into the model by making assumptions on H_1 and H_2 . From the above, we have that $X = (X_1, X_2)$ is a random variable taking values in $(S_1 \times S_2, \mathcal{B}_1 \otimes \mathcal{B}_2)$ (note that in Section 4.3, we consider the example $S_1 = [0, 1]$, $S_2 = [0, 1]^{q-1}$, and $S_1 \times S_2 = [0, 1]^q$, where all spaces are equipped with

the Borel σ -algebra). Moreover, we have that the spaces $L^2(\mathbb{P}^{X_1})$ and $L^2(\mathbb{P}^{X_2})$ are (in a canonical way) subspaces of $L^2(\mathbb{P}^X)$, which implies that H_1 and H_2 are also closed subspaces of $L^2(\mathbb{P}^X)$. Finally, we denote by f the whole regression function given by $f = f_1 + f_2$. We assume that we observe n independent copies

$$(Y^1, X^1), \dots, (Y^n, X^n)$$

of (Y, X) , where $X^i = (X_1^i, X_2^i)$, $1 \leq i \leq n$. Based on this sample, we consider the problem of estimating the function f_1 .

2.2. The main assumption. Our approach relies strongly on the fact that the space $L^2(\mathbb{P}^X)$ is a Hilbert space with the inner product $\langle g, h \rangle = \mathbb{E}[g(X)h(X)]$ and the corresponding norm $\|g\| = \sqrt{\langle g, g \rangle}$ (see, e.g., [9, Theorem 5.2.1]). In order to state our main assumption, we give the following general definition of a minimal angle in Hilbert spaces (see [15, Definition 1] and the references therein).

Definition 1. Let \mathcal{H}_1 and \mathcal{H}_2 be two closed subspaces of a Hilbert space \mathcal{H} with inner product $\langle \cdot, \cdot \rangle$ and norm $\|\cdot\|$. The minimal angle between \mathcal{H}_1 and \mathcal{H}_2 is the number $0 \leq \tau_0 \leq \pi/2$ whose cosine is given by

$$\rho_0 = \rho_0(\mathcal{H}_1, \mathcal{H}_2) = \sup \left\{ \frac{|\langle h_1, h_2 \rangle|}{\|h_1\| \|h_2\|} \mid 0 \neq h_1 \in \mathcal{H}_1, 0 \neq h_2 \in \mathcal{H}_2 \right\}.$$

Assumption 2. Suppose that the cosine of the minimal angle between H_1 and H_2 is strictly less than 1, i.e.,

$$\rho_0(H_1, H_2) < 1.$$

The next lemma states two equivalent formulations of Assumption 2. Since we will also apply it to the finite sample setting in later sections, we again give a general statement.

Lemma 1. Let \mathcal{H}_1 and \mathcal{H}_2 be two closed subspaces of a Hilbert space \mathcal{H} with inner product $\langle \cdot, \cdot \rangle$ and norm $\|\cdot\|$. Let $0 \leq \varrho < 1$ be a constant. Then the following assertions are equivalent:

(i) For all $0 \neq h_1 \in \mathcal{H}_1, 0 \neq h_2 \in \mathcal{H}_2$ we have

$$\frac{|\langle h_1, h_2 \rangle|}{\|h_1\| \|h_2\|} \leq \varrho.$$

(ii) For all $h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2$ we have

$$\|h_1 + h_2\|^2 \geq (1 - \varrho)(\|h_1\|^2 + \|h_2\|^2).$$

(iii) For all $h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2$ we have

$$\|h_1 + h_2\|^2 \geq (1 - \varrho^2)\|h_1\|^2.$$

A proof of Lemma 1 is given in Appendix A.

2.3. The estimation procedure. Let $V_1 \subseteq H_1$ and $V_2 \subseteq H_2$ be d_1 - and d_2 -dimensional linear subspaces, respectively, and let $W_1 \subseteq V_1$ be a linear subspace. Let $V = V_1 + V_2$ and $d = d_1 + d_2$. By Assumption 2, we have $V_1 \cap V_2 = \{0\}$, which implies that d is equal to the dimension of V and that each $g \in V$ can be decomposed uniquely as $g = g_1 + g_2$ with $g_1 \in V_1$ and $g_2 \in V_2$. We will make only one assumption on V which relates the ∞ -norm with the 2-norm, and which will be needed to apply concentration of measure inequalities (compare to, e.g., [5, Section 3.1.1] and [2, Section 1.1]).

Assumption 3. Suppose that there is a real number $\varphi \geq 1$ such that

$$\|g\|_\infty \leq \varphi \sqrt{d} \|g\| \quad (2.1)$$

for all $g \in V$.

Remark 1. In view of Assumption 2, Equation (2.1) is satisfied if there are real numbers $\varphi_j \geq 1$ such that $\|g_j\|_\infty \leq \varphi_j \sqrt{d_j} \|g_j\|$ for all $g_j \in V_j$, $j = 1, 2$. Indeed, applying the Cauchy-Schwarz inequality and Lemma 1, we have

$$\|g_1 + g_2\|_\infty \leq \varphi_1 \sqrt{d_1} \|g_1\| + \varphi_2 \sqrt{d_2} \|g_2\| \leq \frac{\varphi_1 \vee \varphi_2}{\sqrt{1 - \rho_0}} \sqrt{d_1 + d_2} \|g_1 + g_2\|.$$

The construction of our estimator is based on two least squares criteria. First, let \hat{f}_V be the least squares estimator on the model V which is given (not uniquely) by

$$\hat{f}_V = \arg \min_{g \in V} \frac{1}{n} \sum_{i=1}^n (Y^i - g(X^i))^2. \quad (2.2)$$

By the definition of V , we have $\hat{f}_V = (\hat{f}_V)_1 + (\hat{f}_V)_2$ with $(\hat{f}_V)_1 \in V_1$ and $(\hat{f}_V)_2 \in V_2$. Next, by applying a second least squares criterion, we define the estimator \hat{f}_1 by

$$\hat{f}_1 = \arg \min_{g_1 \in W_1} \frac{1}{n} \sum_{i=1}^n ((\hat{f}_V)_1(X_1^i) - g_1(X_1^i))^2. \quad (2.3)$$

We will also consider the special case $W_1 = V_1$, in which we have $\hat{f}_1 = (\hat{f}_V)_1$. This means that the second least squares criterion can be dropped. However, we will see that choosing V_1 as a preliminary space of larger dimension leads to a smaller bias (it lowers the dependence on ρ_0). Finally, since we want to establish risk bounds, it is convenient

to eliminate very large values. Therefore, we define our final estimator \hat{f}_1^* by

$$\hat{f}_1^* = \hat{f}_1 \text{ if } \|\hat{f}_1\|_\infty \leq k_n \text{ and } \hat{f}_1^* \equiv 0 \text{ otherwise,} \quad (2.4)$$

where k_n is a real number to be chosen later (compare to the work by Baraud [2, Eq. (3)]). Finally, note that the estimator is not feasible since the distribution of X is not known and therefore the condition $\mathbb{E}[g_1(X_1)] = 0$ cannot be checked. However, one can replace it by the condition $(1/n) \sum_{i=1}^n g_1(X_1^i) = 0$. In Appendix B, we show how our results carry over to these modified estimators.

In our analysis of \hat{f}_1^* , one important step is to carry over the geometric properties valid in the Hilbert space setting to the finite sample setting. For this, the following event

$$\mathcal{E}_\delta = \{(1 - \delta)\|g\|^2 \leq \|g\|_n^2 \leq (1 + \delta)\|g\|^2 \text{ for all } g \in V\},$$

$0 < \delta < 1$, will play the key role. Here, $\|\cdot\|_n$ denotes the empirical norm (see, e.g., Section 5.1). A first observation is that, under Assumptions 1 and 2, the estimator \hat{f}_1^* is unique on the event \mathcal{E}_δ . This can be seen as follows. If \mathcal{E}_δ holds, then $\|\cdot\|$ and $\|\cdot\|_n$ are equivalent norms on V , which in turn implies that each $g \in V$ is uniquely determined by $(g(X^1), \dots, g(X^n))^T$. Hence, the solutions of the least squares criteria in (2.2) and (2.3) are unique (since the solutions are unique when restricted to vectors in \mathbb{R}^n evaluated at the observations). Moreover, by Assumption 2, the decomposition $\hat{f}_V = (\hat{f}_V)_1 + (\hat{f}_V)_2$ is unique.

In addition, we also obtain a simple representation of our estimator. Let $\hat{\Pi}_V$ be the orthogonal projection from \mathbb{R}^n to the subspace $\{(g(X^1), \dots, g(X^n))^T | g \in V\}$, and let $\hat{\Pi}_{W_1}$ be defined analogously. If \mathcal{E}_δ holds, then we have

$$(\hat{f}_1(X_1^1), \dots, \hat{f}_1(X_1^n))^T = \hat{\Pi}_{W_1}(\hat{\Pi}_V \mathbf{Y})_1,$$

where $\hat{\Pi}_V \mathbf{Y} = (\hat{\Pi}_V \mathbf{Y})_1 + (\hat{\Pi}_V \mathbf{Y})_2$ is the unique decomposition of the least squares estimator on the model V , considered as a vector in \mathbb{R}^n , with $(\hat{\Pi}_V \mathbf{Y})_j \in \{(g_j(X_j^1), \dots, g_j(X_j^n))^T | g_j \in V_j\}$.

3. MAIN RESULTS

3.1. A first risk bound. In this section, we present a first non-asymptotic risk bound in the case $W_1 = V_1$, which will be further improved (under additional assumptions) in later sections. We denote by Π_V (resp. Π_{V_1} , Π_{V_2} , and Π_{W_1}) the orthogonal projection from $L^2(\mathbb{P}^X)$ to the subspace V (resp. V_1 , V_2 , and W_1).

Theorem 1. *Let Assumption 1, 2, and 3 be satisfied. Let $0 < \delta < 1$ be a real number. Let $W_1 = V_1$. Then*

$$\begin{aligned} & \mathbb{E} \left[\|f_1 - \hat{f}_1^*\|^2 \right] \\ & \leq \frac{1+\delta}{(1-\delta)^3} \frac{1}{1-\rho_0^2} \left(\left(1 + \frac{\varphi^2 d}{n} \right) \|f - \Pi_V f\|^2 + \frac{\sigma^2 \dim V_1}{n} \right) + R_n \end{aligned}$$

with

$$R_n =$$

$$\frac{2(1+\delta)\varphi^2 d \|f_1\|^2 (\|f - \Pi_{V_2} f\|^2 + \sigma^2)}{(1-\delta)^2 (1-\rho_0^2) k_n^2} + 2(\|f_1\| + k_n)^2 d \exp \left(-\kappa \frac{\delta^2 n}{\varphi^2 d} \right),$$

where κ is the universal constant in Theorem 7.

Before we discuss the two main terms, let us give conditions under which the remainder term R_n is small. Suppose that for some real number $c > 0$, we have

$$\varphi^2 d \leq \frac{c \delta^2 n}{\log n},$$

and let $k_n^2 = \|f_1\|^2 n^{\kappa/(2c)}$ (this is a theoretical choice of k_n leading to a simple upper bound for R_n , many other choices are possible, too). Then one can show that

$$R_n \leq \frac{12c(1+\delta)\delta^2}{(1-\delta)^2(1-\rho_0^2)} (\|f_1\|^2 + \|f_2 - \Pi_{V_2} f_2\|^2 + \sigma^2) n^{-\frac{\kappa}{2c}+1}.$$

Letting, e.g., $\delta = 1/\log n$ and $c = 1/\log n$, we obtain the following corollary of Theorem 1.

Corollary 1. *Let Assumption 1, 2, and 3 be satisfied. Suppose that*

$$\varphi^2 d \leq \frac{n}{(\log n)^4}. \quad (3.1)$$

Then there is a universal constant $C > 0$ such that

$$\begin{aligned} & \mathbb{E} \left[\|f_1 - \hat{f}_1^*\|^2 \right] \\ & \leq \frac{1}{1-\rho_0^2} \left(\|f_1 - \Pi_{V_1} f_1\|^2 + \frac{\sigma^2 \dim V_1}{n} \right) (1 + C/\log n) \\ & \quad + \frac{C}{1-\rho_0^2} ((\log n) \|f_2 - \Pi_{V_2} f_2\|^2 + \|f_1\|^2 n^{-\frac{\kappa}{2} \log n + 1}). \end{aligned}$$

The first two terms on the right hand side are (up to the factor $(1-\rho_0^2)^{-1}$) equal to the bias term and the variance term of the same estimator with $V_2 = 0$ in the model $Y = f_1(X_1) + \epsilon$. The third term is the approximation error of the function f_2 with respect to the space V_2 .

It decreases if V_2 is chosen larger. Moreover, the choice of V_2 does not effect any of the other terms, the only restriction is given by (3.1). The question arising now is as follows: Is it possible to choose a space V_2 subject to the constraint (3.1) such that $(1 - \rho_0^2)^{-1}(\log n)\|f_2 - \Pi_{V_2}f_2\|^2$ is negligible with respect to the first two terms.

3.2. A refined risk bound. In this section, we improve Theorem 1 such that the factor $(1 - \rho_0^2)^{-1}$ only appears in remainder terms. Since the refined upper bound for the variance term will also contain a Hilbert-Schmidt norm, we give the following general definition (see, e.g., [35]).

Definition 2. Let \mathcal{H}_1 and \mathcal{H}_2 be Hilbert spaces. A bounded linear operator $T : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is called Hilbert-Schmidt if for some orthonormal basis $\{\phi_{1\alpha}\}_{\alpha \in I}$ of \mathcal{H}_1 ,

$$\sum_{\alpha \in I} \|T\phi_{1\alpha}\|^2 < \infty. \quad (3.2)$$

This sum is independent of the choice of the orthonormal basis (see [35, Satz 3.18]). The square root of this sum is called the Hilbert-Schmidt norm of T , denoted by $\|T\|_{HS}$.

Let Π_{V_2} be the orthogonal projection from $L^2(\mathbb{P}^X)$ to V_2 , and let $\Pi_{V_2}|_{W_1}$ be the restriction of Π_{V_2} to W_1 . Then $\Pi_{V_2}|_{W_1}$ is a Hilbert-Schmidt operator, since W_1 is finite-dimensional. We prove:

Theorem 2. *Let Assumption 1, 2, and 3 be satisfied. Let $0 < \delta < 1$ be a real number. Then*

$$\begin{aligned} & \mathbb{E} \left[\|f_1 - \hat{f}_1^*\|^2 \right] \\ & \leq \left(\|f_1 - \Pi_{W_1}f_1\|^2 + \frac{1}{1-\delta} \frac{\sigma^2 \dim W_1}{n} \right) \left(1 + \frac{1+\delta}{(1-\delta)^3} \frac{1}{1-\rho_0^2} \frac{2\varphi^2 d}{n} \right) \\ & + \frac{1+\delta}{(1-\delta)^2} \frac{6}{1-\rho_0^2} (\|f_1 - \Pi_{V_1}f_1\|^2 + \|f_2 - \Pi_{V_2}f_2\|^2) \\ & + \frac{1+\delta}{(1-\delta)^4} \frac{1}{1-\rho_0^2} \frac{\sigma^2 \|\Pi_{V_2}|_{W_1}\|_{HS}^2}{n} + R_n, \end{aligned} \quad (3.3)$$

where R_n is given in Theorem 1.

In order to state a corollary of Theorem 2 similar to Corollary 1, we have to discuss the quantity $\|\Pi_{V_2}|_{W_1}\|_{HS}^2$. If $\{\phi_{1k}\}_{1 \leq k \leq \dim W_1}$ is an orthonormal basis of W_1 , then it can be bounded as follows:

$$\|\Pi_{V_2}|_{W_1}\|_{HS}^2 = \sum_{k=1}^{\dim W_1} \|\Pi_{V_2}\phi_{1k}\|^2 \leq \sum_{k=1}^{\dim W_1} \rho_0^2 \|\phi_{1k}\|^2 = \rho_0^2 \dim W_1, \quad (3.4)$$

where the inequality can be shown as in (5.7). Using this bound, we get a variance term which coincides (up to first order) with the one in Theorem 1. However, (3.4) can be considerably improved under certain Hilbert-Schmidt Assumptions. In particular, we will derive upper bounds which are dimension free. The first assumption is as follows:

Assumption 4. *Suppose that there are measures ν_1 and ν_2 on \mathcal{B}_1 and \mathcal{B}_2 , respectively, such that X has the density p with respect to the product measure $\nu_1 \otimes \nu_2$. Let p_1 and p_2 be the marginal densities of X_1 and X_2 with respect to the measures ν_1 and ν_2 , respectively. Suppose that*

$$\begin{aligned} \|K\|_{HS}^2 &= \int_{S_2} \int_{S_1} \left(\frac{p(x_1, x_2)}{p_1(x_1)p_2(x_2)} \right)^2 p_1(x_1)p_2(x_2) d\nu_1(x_1)d\nu_2(x_2) \\ &= \int_{S_2} \int_{S_1} \frac{(p(x_1, x_2))^2}{p_1(x_1)p_2(x_2)} d\nu_1(x_1)d\nu_2(x_2) < \infty. \end{aligned}$$

If Assumption 4 is satisfied, then we can define the integral operator $K : L^2(\mathbb{P}^{X_1}) \rightarrow L^2(\mathbb{P}^{X_2})$ by

$$(Kg_1)(x_2) = \int_{S_1} g_1(x_1) \frac{p(x_1, x_2)}{p_1(x_1)p_2(x_2)} p_1(x_1) d\nu_1(x_1)$$

which is the orthogonal projection from $L^2(\mathbb{P}^X)$ to $L^2(\mathbb{P}^{X_2})$ restricted to $L^2(\mathbb{P}^{X_1})$. Applying [35, Satz 3.19], we obtain that K is a Hilbert-Schmidt operator with Hilbert-Schmidt norm $\|K\|_{HS}$. We conclude that

$$\|\Pi_{V_2}|_{W_1}\|_{HS} \leq \|K\|_{HS}.$$

Next, we present a more sophisticated upper bound, by using the spaces H_1 and H_2 instead of $L^2(\mathbb{P}^{X_1})$ and $L^2(\mathbb{P}^{X_2})$. Let Π_{H_2} be the orthogonal projection from $L^2(\mathbb{P}^X)$ to H_2 , and let $\Pi_{H_2}|_{H_1}$ be the restriction of Π_{H_2} to H_1 .

Assumption 5 (Weaker form of Assumption 4). *Suppose that $\Pi_{H_2}|_{H_1}$ is a Hilbert-Schmidt operator.*

If Assumption 5 is satisfied, then

$$\|\Pi_{V_2}|_{W_1}\|_{HS} \leq \|\Pi_{H_2}|_{H_1}\|_{HS}.$$

Letting now $\delta = 1/\log n$ and $c = 1/\log n$ as in Corollary 1, we obtain the following corollary of Theorem 2.

Corollary 2. *Let Assumption 1, 2, 3, and 4 be satisfied. Suppose that*

$$\varphi^2 d \leq \frac{n}{(\log n)^4}.$$

Then there is a universal constant $C > 0$ such that

$$\begin{aligned} & \mathbb{E} \left[\|f_1 - \hat{f}_1^*\|^2 \right] \\ & \leq \left(\|f_1 - \Pi_{W_1} f_1\|^2 + \frac{\sigma^2 \dim W_1}{n} \right) (1 + C' / \log n) \\ & \quad + C' \left(\|f_1 - \Pi_{V_1} f_1\|^2 + \|f_2 - \Pi_{V_2} f_2\|^2 + \frac{\sigma^2 \|K\|_{HS}^2}{n} + \frac{\|f_1\|^2}{n^{\frac{\kappa}{2} \log n - 1}} \right), \end{aligned}$$

where $C' = C/(1 - \rho_0^2)$. Moreover, if Assumption 5 holds instead of Assumption 4, then the above inequality holds if $\|K\|_{HS}^2$ is replaced by $\|\Pi_{H_2}|_{H_1}\|_{HS}^2$. Finally, if Assumption 5 and 4 are not satisfied, then the above inequality holds if $\|K\|_{HS}^2$ is replaced by $\rho_0^2 \dim W_1$.

Now the first two terms in the brackets on the right hand side are equal to the bias term and the variance term of the same estimator with $V_2 = 0$ in the model $Y = f_1(X_1) + \epsilon$. As in Corollary 1, we see that the choices of V_1 and V_2 do not effect any of the other terms, the only restriction is given by (3.1).

Finally, we give an alternative representation of the Hilbert-Schmidt norm $\|\Pi_{H_2}|_{H_1}\|_{HS}$ using the operator $\Pi_{H_1}\Pi_{H_2}\Pi_{H_1}$ (which we consider as a map from H_1 to H_1). To simplify the exposition, we suppose that H_1 is separable, which implies that each orthonormal basis of H_1 is countable (see, e.g., [27, Chapter II]). From Assumption 5, it follows that $\Pi_{H_1}\Pi_{H_2}\Pi_{H_1}$ is compact (see, e.g., [18, Chapter 30.8]). Since it is also symmetric and positive, the spectral theorem (see, e.g., [18, Theorem 3 in Chapter 28]) implies that there is an orthonormal basis for H_1 consisting of eigenvectors of $\Pi_{H_1}\Pi_{H_2}\Pi_{H_1}$. These all have non-negative eigenvalues. We arrange the positive eigenvalues of $\Pi_{H_1}\Pi_{H_2}\Pi_{H_1}$ in decreasing order: $\alpha_1 \geq \alpha_2 \cdots > 0$. We now have:

Lemma 2. *Under the above assumptions, we have*

$$\rho_0^2 = \alpha_1$$

and

$$\|\Pi_{H_2}|_{H_1}\|_{HS}^2 = \text{tr}(\Pi_{H_1}\Pi_{H_2}\Pi_{H_1}) = \sum_{k \geq 1} \alpha_k.$$

Proof. We only prove the second equality. Let $\{\phi_{1k}\}_{k \geq 1}$ be an orthonormal basis for H_1 consisting of eigenvectors of $\Pi_{H_1}\Pi_{H_2}\Pi_{H_1}$. Then

$$\|\Pi_{H_2}|_{H_1}\|_{HS}^2 = \sum_{k \geq 1} \|\Pi_{H_2}\phi_{1k}\|^2 = \sum_{k \geq 1} \langle \Pi_{H_1}\Pi_{H_2}\Pi_{H_1}\phi_{1k}, \phi_{1k} \rangle^2 = \sum_{k \geq 1} \alpha_k.$$

□

Example 1. Consider the case that $X = (X_1, X_2)$ is a bivariate Gaussian random variable such that $\mathbb{E}[X_1] = \mathbb{E}[X_2] = 0$, $\mathbb{E}[X_1^2] = \mathbb{E}[X_2^2] = 1$, and $\mathbb{E}[X_1 X_2] = \rho$.

First, suppose that H_1 and H_2 are the spaces of linear centered functions, i.e., $H_1 = \{g_1 : g_1(x_1) = a \cdot x_1, a \in \mathbb{R}\}$ and $H_2 = \{g_2 : g_2(x_2) = a \cdot x_2, a \in \mathbb{R}\}$. Then it is easy to see that

$$\rho_0 = |\rho|$$

and

$$\|\Pi_{H_2}|_{H_1}\|_{HS}^2 = \rho^2.$$

Second, suppose that $H_1 = \{g_1 \in L^2(\mathbb{P}^{X_1}) : \mathbb{E}[g_1(X_1)] = 0\}$ and $H_2 = L^2(\mathbb{P}^{X_2})$. Then it follows from [17] that $\Pi_{H_1}\Pi_{H_2}\Pi_{H_1}$ has eigenvalues $\{\rho^2, \rho^4, \dots\}$. Hence, the above lemma implies that

$$\rho_0 = |\rho|$$

and

$$\|\Pi_{H_2}|_{H_1}\|_{HS}^2 = \sum_{k=1}^{\infty} \rho^{2k} = \frac{\rho^2}{1 - \rho^2},$$

which is an improvement over (3.4) if $\dim W_1$ is large.

3.3. Regularity conditions on the design densities. In this section, we present two improvements of Theorem 2 which are possible under Assumption 4 and additional regularity conditions on the design densities. In particular, we show that the dependence of the bias term on the function f_2 can decrease considerably.

By Assumption 4 and Fubini's theorem, we have

$$\frac{p(x_1, \cdot)}{p_1(x_1)p_2(\cdot)} \in L^2(\mathbb{P}^{X_2}) \quad (3.5)$$

for \mathbb{P}^{X_1} -almost all x_1 . Thus we can make the following assumption. Suppose that there is a real number $\psi(V_2)$ and a function $h_1 \in L^2(\mathbb{P}^{X_1})$ such that

$$\left\| (1 - \Pi_{V_2}) \frac{p(x_1, \cdot)}{p_1(x_1)p_2(\cdot)} \right\|_{L^2(\mathbb{P}^{X_2})} \leq h_1(x_1)\psi(V_2) \quad (3.6)$$

for \mathbb{P}^{X_1} -almost all x_1 . In analogy, we let $\phi(V_2)$ be a real number such that $\|f_2 - \Pi_{V_2}f_2\| \leq \phi(V_2)$. We prove:

Theorem 3. *Let Assumption 1, 2, 3, and 4 be satisfied. Let $0 < \delta < 1$ be a real number. Suppose that (3.6) is satisfied. Moreover, suppose*

that $\|g_1\|_\infty \leq \varphi\sqrt{d_1}\|g_1\|$ for all $g_1 \in V_1$, where φ is the constant from Assumption 3. Then (3.3) holds when

$$\frac{1+\delta}{(1-\delta)^2} \frac{6}{1-\rho_0^2} (\|f_1 - \Pi_{V_1} f_1\|^2 + \|f_2 - \Pi_{V_2} f_2\|^2)$$

is replaced by

$$\begin{aligned} & \frac{(1+\delta)^2}{(1-\delta)^4} \frac{12}{1-\rho_0^2} \left(\|f_1 - \Pi_{V_1} f_1\|^2 + \frac{\|h_1\|^2 (\phi(V_2)\psi(V_2))^2}{1-\rho_0^2} \right. \\ & \quad \left. + \frac{1}{n} \frac{\|h_1\|^2 \|f_2 - \Pi_{V_2} f_2\|_\infty^2 (\psi(V_2))^2}{1-\rho_0^2} + \frac{(\phi(V_2))^2}{1-\rho_0^2} \frac{\varphi^2 d_1}{n} \right). \end{aligned}$$

Theorem 3 shows that the regularity conditions on $p/(p_1 p_2)$ and f_2 have similar effects, which can be seen from second term. In contrast to Theorems 1 and 2, Theorem 3 shows that the estimator \hat{f}_1^* can also behave well when f_2 is considerably less regular than f_1 . For instance, if we apply Theorem 3 to an asymptotic scenario, then, under suitable conditions on $\psi(V_2)$, the regularity conditions on f_2 can be (almost) reduced to $\phi(V_2) \rightarrow 0$ (see, e.g., Corollary 6).

For fixed x_1 , let the function $r(x_1, \cdot)$ be the orthogonal projection of $p(x_1, \cdot)/(p_1(x_1)p_2(\cdot))$ on H_2 . By (3.5), $r(x_1, \cdot)$ is defined for \mathbb{P}^{X_1} -almost all x_1 . Thus we can consider the following weaker version of (3.6). Suppose that there exists a real number $\psi_\Pi(V_2)$ and a function $h_1 \in L^2(\mathbb{P}^{X_1})$ such that

$$\|(1 - \Pi_{V_2})r(x_1, \cdot)\|_{L^2(\mathbb{P}^{X_2})} \leq h_1(x_1)\psi_\Pi(V_2) \quad (3.7)$$

for \mathbb{P}^{X_1} -almost all x_1 . If (3.7) holds, then we obtain the following theorem. Note that, compared to Theorem 3, the last term is not always negligible.

Theorem 4. *Let Assumption 1, 2, 3, and 4 be satisfied. Let $0 < \delta < 1$ be a real number. Suppose that (3.7) is satisfied. Then Theorem 2 also holds when the term*

$$\frac{1+\delta}{(1-\delta)^2} \frac{6}{1-\rho_0^2} (\|f_1 - \Pi_{V_1} f_1\|^2 + \|f_2 - \Pi_{V_2} f_2\|^2)$$

is replaced by

$$\begin{aligned} & \frac{1+\delta}{(1-\delta)^3} \frac{6}{1-\rho_0^2} \left(\|f_1 - \Pi_{V_1} f_1\|^2 \left(1 + \frac{2\varphi^2 d}{n} \right) \right. \\ & \quad \left. + \frac{\|h_1\|^2 (\phi(V_2)\psi_\Pi(V_2))^2}{1-\rho_0^2} + \frac{(\phi(V_2))^2}{1-\rho_0^2} \frac{2\varphi^2 d}{n} \right). \end{aligned}$$

4. APPLICATIONS

4.1. The two-dimensional case. In this section, we want to discuss Theorem 1 and 2 in the case that X_1 and X_2 take values in \mathbb{R} and that Assumptions 1 and 2 are satisfied with

$$H_1 = \{g_1 \in L^2(\mathbb{P}^{X_1}) \mid \mathbb{E}[g_1(X_1)] = 0\}$$

and

$$H_2 = L^2(\mathbb{P}^{X_2}).$$

The main remaining issue is to bound the approximation errors. This is possible if the f_j belong to certain nonparametric classes of functions and if the V_j are chosen appropriately. Here, we shall restrict our attention to (periodic) Sobolev smoothness and spaces of trigonometric polynomials. Note that we will also consider Hölder smoothness and spaces of piecewise polynomials in Section 4.4 and 4.5. Recall that the trigonometric basis is given by $\phi_0(x) = 1$, $\phi_k(x) = \sqrt{2} \cos(2\pi kx)$ and $\phi_{-k}(x) = \sqrt{2} \sin(2\pi kx)$, $k \geq 1$, where $x \in [0, 1]$.

Assumption 6. *Suppose that the X_j take values in $[0, 1]$ and have densities p_{X_j} with respect to the Lebesgue measure on $[0, 1]$, which satisfy $c \leq p_{X_j} \leq 1/c$ for some constant $c > 0$. Moreover, suppose that the f_j belong to the Sobolev classes*

$$\tilde{W}_j(\alpha_j, K_j) = \left\{ \sum_{k \in \mathbb{Z}} \theta_k \phi_k(x_j) : \sum_{k \in \mathbb{Z}} |k|^{2\alpha_j} \theta_k^2 \leq K_j^2 \right\},$$

where $\alpha_j > 0$ and $K_j > 0$ (see, e.g., [32, Definition 1.12]).

For $j = 1, 2$, let V_j be the intersection of H_j with the linear span of the ϕ_k (in the variable x_j) such that $|k| \leq m_j$, and let W_1 be the intersection of H_1 with the linear span of the ϕ_k (in the variable x_1) such that $|k| \leq m_{W_1}$. Note that we have $d_1 = 2m_1$ and $d_2 = 2m_2 + 1$. Using the definition of the $\tilde{W}_j(\alpha_j, K_j)$, we have for all $h_j \in \tilde{W}_j(\alpha_j, K_j) \cap H_j$,

$$\|h_j - \Pi_{V_j} h_j\|^2 \leq (1/c) K_j^2 (1 + m_j)^{-2\alpha_j}$$

and thus

$$\|h_j - \Pi_{V_j} h_j\|^2 \leq C_j K_j^2 d_j^{-2\alpha_j}, \quad (4.1)$$

where $C_j = 2^{2\alpha_j}/c$. The same bound holds if V_1 and d_1 are replaced by W_1 and $\dim W_1$. Moreover, by applying the Cauchy-Schwarz inequality, we have $\|g_j\|_\infty \leq \sqrt{1/c} \sqrt{2m_j + 1} \|g_j\|$ for all $g_j \in V_j$, which implies that $\|g_j\|_\infty \sqrt{2/c} \sqrt{d_j} \|g_j\|$ for all $g_j \in V_j$. By Remark 1, this implies that Assumption 3 is satisfied with

$$\varphi^2 \leq \frac{2}{c(1 - \rho_0)}. \quad (4.2)$$

We now choose d_1 and d_2 as the smallest possible integer satisfying

$$d_j \geq \frac{n}{4\varphi^2 \log^4 n}, \quad (4.3)$$

and we choose $\dim W_1$ (up to constant) equal to

$$\left(\frac{K_1^2 n}{\sigma^2} \right)^{\frac{1}{2\alpha_1+1}}.$$

If n is large enough, then these choices imply that (3.1) of Corollary 1 is satisfied. Applying (4.1) and (4.3), we obtain that

$$\|f_2 - \Pi_{V_2} f_2\|^2 \leq C_2 K_2^2 \left(\frac{n}{4\varphi^2 \log^4 n} \right)^{-2\alpha_2},$$

where the last expression is $o(n^{-(\alpha_1/(2\alpha_1+1))})$ if $\alpha_2 > \alpha_1/(2\alpha_1+1)$. Finally, note that $\|f_1\|^2 \leq \mathbb{E}[(f_1(X_1) - \theta_0)^2] \leq K_1^2/c$. From Corollary 2, we now obtain the following asymptotic result when the sample size n tends to infinity:

Corollary 3. *Let Assumption 2 and 6 be satisfied. Suppose that*

$$\alpha_2 > \alpha_1/(2\alpha_1+1). \quad (4.4)$$

Then

$$\limsup_{n \rightarrow \infty} \sup_{f_j \in \tilde{W}_j(\alpha_j, K_j)} \mathbb{E} \left[n^{\frac{2\alpha_1}{2\alpha_1+1}} \|f_1 - \hat{f}_1^*\|^2 \right] \leq C \sigma^{\frac{4\alpha_1}{2\alpha_1+1}} K_1^{\frac{2}{2\alpha_1+1}}, \quad (4.5)$$

where C is a constant depending only on α_1 , c , and ρ_0 . If, in addition, Assumption 5 holds, then the dependence of C on ρ_0 disappears and we obtain the same constant as for the corresponding estimator with $V_2 = 0$ in the model $Y = f_1(X_1) + \epsilon$.

Remark 2. Corollary 3 says that the estimator \hat{f}_1^* attains the optimal rate of convergence in a minimax sense. Note that one can also take the supremum over random variables X such that c , $1/c$, and $1/(1-\rho_0^2)$ are bounded by a fixed constant.

Remark 3. The assumptions of Corollary 3 are weaker than those needed in [13], where it is also assumed that the joint density of (X_1, X_2) is bounded away from zero and infinity and that the functions f_1 and f_2 are Lipschitz continuous. Note that (4.4) is always satisfied if, e.g., $\alpha_2 \geq 1/2$ and that the boundedness conditions imply Assumption 4 and also Lemma 1 (ii). Hence, the boundedness conditions imply Assumption 2 and 5.

Remark 4. In semiparametric models, one often requires a global rate of convergence of order $o(n^{-1/4})$ in order to obtain the rate of convergence $n^{-1/2}$ for the parametric component (see, e.g., the book by van de Geer [33, Chapter 11]). Since $\alpha_1/(2\alpha_1 + 1)$ goes to $1/2$ as $\alpha_1 \rightarrow \infty$, condition (4.4) extends this result to the nonparametric case.

4.2. The multidimensional case. Now, we suppose that the X_1 and X_2 take values in $[0, 1]^{q_1}$ and $[0, 1]^{q_2}$, respectively, and that Assumptions 1 and 2 are again satisfied with

$$H_1 = \{g_1 \in L^2(\mathbb{P}^{X_1}) | \mathbb{E}[g_1(X_1)] = 0\}$$

and

$$H_2 = L^2(\mathbb{P}^{X_2}).$$

In this case, we consider the tensor product Fourier basis:

$$\phi_k(x_j) = \prod_{l=1}^{q_j} \phi_{k_l}(x_{jl}),$$

where $k \in \mathbb{Z}^{q_j}$ and $x_j \in [0, 1]^{q_j}$. We define the following Sobolev class

$$\tilde{W}_j(\alpha_j, K_j) = \left\{ \sum_{k \in \mathbb{Z}^{q_j}} \theta_k \phi_k(x_j) : \sum_{k \in \mathbb{Z}^{q_j}} a_{jk}^2 \theta_k^2 \leq K_j^2 \right\}$$

with $a_{jk} = \|k\|_\infty^{\alpha_j}$, where $\|k\|_\infty = \max_{l=1, \dots, q_j} |k_l|$, $\alpha_j > 0$, and $K_j > 0$ (an alternative choice would be, e.g., $a_{jk}^2 = \sum_{l=1}^{q_j} |k_l|^{\alpha_j}$).

Assumption 7. Suppose that the X_j have densities p_{X_j} with respect to the Lebesgue measure on $[0, 1]^{q_j}$, which satisfy $c \leq p_{X_j} \leq 1/c$ for some constant $c > 0$. Moreover, suppose that the f_j belong to the Sobolev classes $\tilde{W}_j(\alpha_j, K_j)$, where $\alpha_j > 0$ and $K_j > 0$.

For $j = 1, 2$, let V_j be the intersection of H_j with the linear span of the ϕ_k (in the variable x_j) such that $\|k\|_\infty \leq m_j$, and let W_1 be the intersection of H_1 with the linear span of the ϕ_k (in the variable x_1) such that $\|k\|_\infty \leq m_{W_1}$. Note that we have $d_1 = (2m_1 + 1)^{q_1} - 1$ and $d_2 = (2m_2 + 1)^{q_2}$. Similarly as above, one can show that

$$\|h_j - \Pi_{V_j} h_j\|^2 \leq C_j K_j^2 d_j^{-2\alpha_j/q_j}$$

for all $h_j \in \tilde{W}_j(\alpha_j, K_j) \cap H_j$, where $C_j = 2^{2\alpha_j}/c$, and that Assumption 3 holds with a φ satisfying again (4.2). We now choose d_1 and d_2 as in (4.3), and we choose $\dim W_1$ (up to constant) equal to

$$\left(\frac{K_1^2 n}{\sigma^2} \right)^{\frac{q_1}{2\alpha_1 + q_1}}.$$

Similarly as above, we conclude:

Corollary 4. *Let Assumption 2 and 7 be satisfied. Suppose that*

$$\alpha_2/q_2 > \alpha_1/(2\alpha_1 + q_1).$$

Then

$$\limsup_{n \rightarrow \infty} \sup_{f_j \in \tilde{W}_j(\alpha_j, K_j)} \mathbb{E} \left[n^{\frac{2\alpha_1}{2\alpha_1 + q_1}} \|f_1 - \hat{f}_1^*\|^2 \right] \leq C \sigma^{\frac{4\alpha_1}{2\alpha_1 + q_1}} K_1^{\frac{2q_1}{2\alpha_1 + q_1}},$$

where C is a constant depending only on α_1 , c , and ρ_0 . If, in addition, Assumption 5 holds, then the dependence of C on ρ_0 disappears and we obtain the same constant as for the corresponding estimator with $V_2 = 0$ in the model $Y = f_1(X_1) + \epsilon$.

4.3. The additive model with Sobolev smoothness. In this section, we want to discuss Theorem 1 and 2 in the case that the random variables X_1 and X_2 take values in \mathbb{R} and \mathbb{R}^{q-1} , $q \geq 2$, respectively, and that Assumptions 1 and 2 are satisfied with

$$H_1 = \{g_1 \in L^2(\mathbb{P}^{X_1}) | \mathbb{E}[g_1(X_1)] = 0\}$$

and

$$H_2 = \sum_{j=1}^{q-1} L^2(\mathbb{P}^{X_{2j}}),$$

where the X_{2j} , $j = 1, \dots, q-1$, are the components of X_2 . If we define $H_{2j} = \{g_{2j} \in L^2(\mathbb{P}^{X_{2j}}) | \mathbb{E}[g_{2j}(X_{2j})] = 0\}$, if $j = 1, \dots, q-2$, and $H_{2j} = L^2(\mathbb{P}^{X_{2j}})$, if $j = q-1$, then we can write $H_2 = \sum_{j=1}^{q-1} H_{2j}$.

Assumption 8. *Suppose that X_1 and the X_{2j} take values in $[0, 1]$ and have densities p_{X_1} and $p_{X_{2j}}$ with respect to the Lebesgue measure on $[0, 1]$, which satisfy $c \leq p_{X_1} \leq 1/c$ and $c \leq p_{X_{2j}} \leq 1/c$ for some constant $c > 0$.*

Moreover, suppose that $f_1 \in \tilde{W}_1(\alpha_1, K_1)$, where $\alpha_1 > 0$ and $K_1 > 0$, and that there is a decomposition $f_2 = \sum_{j=1}^{q-1} f_{2j}$ such that $f_{2j} \in \tilde{W}_{2j}(\alpha_2, K_2) \cap H_{2j}$, where $\alpha_2 > 0$ and $K_2 > 0$.

Now, let V_1 and W_1 be as in Section 4.1, and let $V_2 = \sum_{j=1}^{q-1} V_{2j}$, where V_{2j} is the intersection of H_{2j} with the linear span of the ϕ_k (in the variable x_{2j}) such that $|k| \leq m_2$. In order to show that Assumption 3 is satisfied, we need the following additional assumption on H_2 .

Assumption 9. *There is an $\epsilon_2 < 1$ such that for each $h_2 \in H_2$, there is a decomposition $h_2 = \sum_{j=1}^{q-1} h_{2j}$ with $h_{2j} \in H_{2j}$ such that*

$$\|h_2\|^2 \geq (1 - \epsilon_2) \sum_{j=1}^{q-1} \|h_{2j}\|^2. \quad (4.6)$$

By applying iteratively [3, Proposition 2.A in Appendix A.4], one can show that Assumption 9 is equivalent to the assertion that $\sum_{j \in J} H_{2j}$ is closed for all $J \subseteq \{1, \dots, q-1\}$. In particular, Assumption 9 implies that H_2 is closed meaning that Assumption 1 is included in Assumption 9. We mention that Assumption 9 can also be related to bounds on certain complementary angles (see [3, Definition 2 and Proposition 2.D in Appendix A.4]).

Lemma 3. *Let Assumption 2, 8, and 9 be satisfied. Then Assumption 3 holds with a constant φ satisfying*

$$\varphi^2 \leq \frac{2}{c(1-\rho_0)(1-\epsilon_2)} \frac{d_1 + \sum_{j=1}^{q-1} d_{2j}}{d},$$

where $d_{2j} = \dim V_{2j}$. In particular, if the V_{2j} are linearly independent, then we have

$$\varphi^2 \leq \frac{2}{c(1-\rho_0)(1-\epsilon_2)}. \quad (4.7)$$

A proof of Lemma 3 is given in Appendix C. Thus (3.1) of Corollary 1 is satisfied if

$$\frac{2 \left(d_1 + \sum_{j=1}^{q-1} d_{2j} \right)}{c(1-\rho_0)(1-\epsilon_2)} \leq \frac{n}{\log^4 n}.$$

We now choose $\dim W_1$ as in Section 4.1, and we choose d_1 and the d_{2j} equal to the smallest possible integers satisfying

$$d_1 \geq \frac{c(1-\rho_0)(1-\epsilon_2)n}{8 \log^4 n}$$

and

$$d_{2j} \geq \frac{c(1-\rho_0)(1-\epsilon_2)n}{8(q-1) \log^4 n}. \quad (4.8)$$

If the right hand side of (4.8) is greater than or equal to 2, then (3.1) of Corollary 1 is satisfied. In order to bound the approximation error $\|f_2 - \Pi_{V_2} f_2\|$, we introduce ϵ'_2 which is the smallest number such that

$$\|h_2\|^2 \leq (1 + \epsilon'_2) \sum_{j=1}^{q-1} \|h_{2j}\|^2 \quad (4.9)$$

for all $h_2 = \sum_{j=1}^{q-1} h_{2j}$ with $h_{2j} \in H_{2j}$. Note that

$$1 + \epsilon'_2 \leq q - 1,$$

by the Cauchy-Schwarz inequality. Using the decomposition of f_2 in Assumption 8, the projection theorem, (4.9), and finally (4.1) and (4.8),

we have

$$\begin{aligned}
\|f_2 - \Pi_{V_2} f_2\|^2 &\leq \left\| \sum_{j=1}^{q-1} (f_{2j} - \Pi_{V_{2j}} f_{2j}) \right\|^2 \\
&\leq (1 + \epsilon'_2) \sum_{j=1}^{q-1} \|f_{2j} - \Pi_{V_{2j}} f_{2j}\|^2 \\
&\leq (1 + \epsilon'_2) C_2 K_2^2 (q-1) \left(\frac{c(1 - \rho_0)(1 - \epsilon_2)n}{8(q-1) \log^4 n} \right)^{-2\alpha_2}.
\end{aligned}$$

From Corollary 2, we now obtain:

Theorem 5. *Let Assumption 2, 8, and 9 be satisfied. Suppose that the right hand side of (4.8) is greater than or equal to 2 and that $\dim W_1 \leq d_1$. Then*

$$\begin{aligned}
\sup_{f_1 \in \tilde{W}_1(\alpha_1, K_1)} \sup_{f_{2j} \in \tilde{W}_{2j}(\alpha_2, K_2)} \mathbb{E} \left[\|f_1 - \hat{f}_1^*\|^2 \right] &\leq C \sigma^{\frac{4\alpha_1}{2\alpha_1+1}} K_1^{\frac{2}{2\alpha_1+1}} n^{-\frac{2\alpha_1}{2\alpha_1+1}} \\
&\quad + C'(1 + \epsilon'_2) q^{2\alpha_2+1} (\log n)^{8\alpha_2} ((1 - \epsilon_2)n)^{-2\alpha_2}, \quad (4.10)
\end{aligned}$$

where C is a constant depending only on α_1 , c and ρ_0 , and C' is a constant depending only on α_2 , c , ρ_0 , K_2 , and K_1 .

Remark 5. If the last expression in (4.10) is of smaller order, then Theorem 5 says that the estimator \hat{f}_1^* attains the (nonasymptotic) optimal rate of convergence in a minimax sense. The last expression in (4.10) is of smaller order if, e.g.,

$$q^{2\alpha_2+1} (\log n)^{8\alpha_2+1} n^{-2\alpha_2} \leq C'' n^{-\frac{2\alpha_1}{2\alpha_1+1}}, \quad (4.11)$$

where C'' depends only on α_2 , c , ρ_0 , K_2 , K_1 , ϵ_2 , and ϵ'_2 . This result can be applied to an asymptotic scenario in which q and n tend to infinity such that (4.11) is satisfied.

Next, we apply Theorem 5 in the case that the sample size n tends to infinity, and q is a fixed constant.

Corollary 5. *Let Assumption 2, 8, and 9 be satisfied. Suppose that*

$$\alpha_2 > \alpha_1 / (2\alpha_1 + 1).$$

Then

$$\limsup_{n \rightarrow \infty} \sup_{f_1 \in \tilde{W}_1(\alpha_1, K_1)} \sup_{f_{2j} \in \tilde{W}_{2j}(\alpha_2, K_2)} \mathbb{E} \left[n^{\frac{2\alpha_1}{2\alpha_1+1}} \|f_1 - \hat{f}_1^*\|^2 \right] \leq C \sigma^{\frac{4\alpha_1}{2\alpha_1+1}} K_1^{\frac{2}{2\alpha_1+1}},$$

where C is a constant depending only on α_1 , c , and ρ_0 . If, in addition, Assumption 5 holds, then the dependence of C on ρ_0 disappears and

we obtain the same constant as for the corresponding estimator with $V_2 = 0$ in the model $Y = f_1(X_1) + \epsilon$.

Remark 6. Again, the assumptions of Corollary 5 are weaker than those needed in [13] (compare to Remark 3). For instance, the condition that the one- and two-dimensional marginal densities of X_{2j} and $(X_{2j}, X_{2j'})$ are bounded away from zero and infinity implies that Assumption 9 is satisfied (see, e.g., [3, Proposition 2.C in Appendix A.4] for the case $q = 3$ and [20, Lemma 1] for the general case).

4.4. The additive model with Hölder smoothness. We continue the discussion of the additive model in Section 4.3, and consider briefly the case of Hölder smoothness and spaces of piecewise polynomials.

Assumption 10. Suppose that X_1 and the X_{2j} take values in $[0, 1]$ and have densities p_{X_1} and $p_{X_{2j}}$ with respect to the Lebesgue measure on $[0, 1]$, which satisfy $p_{X_1} \geq c$ and $p_{X_{2j}} \geq c$ for some constant $c > 0$.

Moreover, suppose that the function f_1 belongs to the Hölder class $\mathcal{H}_1(\alpha_1, K_1)$ on $[0, 1]$, where $\alpha_1 > 0$ and $K_1 \geq 0$ (see, e.g., [32, Definition 1.2]), and that there is a decomposition $f_2 = \sum_{j=1}^{q-1} f_{2j}$ such that $f_{2j} \in \mathcal{H}_{2j}(\alpha_2, K_2) \cap H_{2j}$, where $\alpha_2 > 0$ and $K_2 \geq 0$.

Let V_1 be the intersection of H_1 with the space of regular piecewise polynomials (in the variable x_1) with integer-valued parameters $r_1 = \lfloor \alpha_1 \rfloor$ and m_1 , where r_1 is the maximal degree of the polynomials and $\{0 < 1/m_1 < 2/m_1 < \dots < 1\}$ generates the partition of $[0, 1]$ into m_1 intervals (see, e.g., [4]), and let W_1 be the intersection of H_1 with the space of regular piecewise polynomials (in the variable x_1) with integer-valued parameters $r_1 = \lfloor \alpha_1 \rfloor$ and m_{W_1} (in order that $W_1 \subseteq V_1$ we need that m_1 is a multiple of m_{W_1}). Moreover, let $V_2 = \sum_{j=1}^{q-1} V_{2j}$, where V_{2j} is the intersection of H_{2j} with the space of regular piecewise polynomials (in the variable x_{2j}) with integer-valued parameters $r_2 = \lfloor \alpha_2 \rfloor$ and m_2 . Note that alternatively, one could also consider spaces of splines with the same parameters (see, e.g., [8, Chapter VII, VIII]). We have $d_1 = (r_1 + 1)m_1 - 1$, $d_{2j} = (r_2 + 1)m_2 - 1$, if $j = 1, \dots, q-1$, and $d_{2j} = (r_2 + 1)m_2$, if $j = q-1$. Using Taylor's theorem, one can show that there are constants C_1 and C_2 depending only on α_1 and α_2 , respectively, such that

$$\inf_{g_1 \in V_1} \|h_1 - g_1\|_\infty^2 \leq C_1 K_1^2 d_1^{-2\alpha_1} \quad (4.12)$$

for all $h_1 \in \mathcal{H}(\alpha_1, K_1) \cap H_1$ and

$$\inf_{g_{2j} \in V_{2j}} \|h_{2j} - g_{2j}\|_\infty^2 \leq C_2 K_2^2 d_{2j}^{-2\alpha_2} \quad (4.13)$$

for all $h_{2j} \in \mathcal{H}(\alpha_2, K_2) \cap H_{2j}$. Note that similar bounds also hold for spline spaces (see, e.g., [8, Chapter XII]). Concerning Assumption 3, we have a similar result as in the previous section:

Lemma 4. *Let Assumptions 2 and 9, and 10 be satisfied. Let $r = r_1 \vee r_2$. Then we have*

$$\varphi^2 \leq \frac{2(r+1)}{c(1-\rho_0)(1-\epsilon_2)} \frac{d_1 + \sum_{j=1}^{q-1} d_{2j}}{d}.$$

In particular, if the V_{2j} are linearly independent, then we have

$$\varphi^2 \leq \frac{2(r+1)}{c(1-\rho_0)(1-\epsilon_2)}.$$

A proof of Lemma 4 is given in Appendix C. Choosing d_1 , $\dim W_1$, and d_{2j} , $j = 1, \dots, q-1$, similarly as in the previous section, we obtain the following analogue of Theorem 5:

Theorem 6. *Let Assumption 2, 9, and 10 be satisfied. Suppose that the right hand side of (4.8) is greater than or equal to $r+1$ and that $\dim W_1 \leq c(1-\rho_0)(1-\epsilon_2)n/(4(r+1)\log^4 n)$. Then*

$$\begin{aligned} \mathbb{E} \left[\|f_1 - \hat{f}_1^*\|^2 \right] &\leq C \sigma^{\frac{4\alpha_1}{2\alpha_1+1}} K_1^{\frac{2}{2\alpha_1+1}} n^{-\frac{2\alpha_1}{2\alpha_1+1}} \\ &\quad + C'(1+\epsilon'_2) \varphi^{4\alpha_2} q^{2\alpha_2+1} (\log n)^{8\alpha_2+1} ((1-\epsilon_2)n)^{-2\alpha_2} \end{aligned}$$

where C is a constant depending only on α_1 and ρ_0 , and C' is a constant depending only on α_2 , c , r , ρ_0 , K_2 , and $\|f_1\|$.

4.5. The additive model with smooth design densities. We continue the example in Section 4.4 and discuss Condition (3.6) and (3.7) in Theorem 3 and 4, respectively. First, we apply Theorem 3 in the simple case $q = 2$. We suppose that Assumption 4 holds and that for each fixed x_1 ,

$$\frac{p_X(x_1, x_2)}{p_{X_1}(x_1)p_{X_2}(x_2)} \in \mathcal{H}(\beta, h_1(x_1)) \quad (4.14)$$

with $h_1 \in L^2(\mathbb{P}^{X_1})$. Let V_1 and W_1 as in the previous section, and let V_2 be the space of regular piecewise polynomials in the variable x_2 with parameters $r_2 = \lfloor \alpha_2 \rfloor \vee \lfloor \beta \rfloor$ and d_2 as in (4.3). Then (4.13) holds with a constant C_2 depending only on α_2 and r_2 . Moreover, (3.6) is satisfied with h_1 from (4.14) and $\psi(V_2) = \sqrt{C_3} d_2^{-\beta}$, where C_3 is a constant depending only on β and r_2 . Thus

$$\|h_1\| \psi(V_2) \phi(V_2) \leq \sqrt{C_2 C_3} K_2 \|h_1\| \left(\frac{n}{4\varphi^2 \log^4 n} \right)^{-\alpha_2 - \beta}$$

From Theorem 3, we obtain:

Corollary 6. *Let $q = 2$. Let Assumption 2, 10, and 4 be satisfied. Suppose that (4.14) holds, and that*

$$\alpha_2 + \beta > \alpha_1/(2\alpha_1 + 1). \quad (4.15)$$

Then

$$\limsup_{n \rightarrow \infty} \mathbb{E} \left[n^{\frac{2\alpha_1}{2\alpha_1+1}} \|f_1 - \hat{f}_1^*\|^2 \right] \leq C \sigma^{\frac{4\alpha_1}{2\alpha_1+1}} K_1^{\frac{2}{2\alpha_1+1}}, \quad (4.16)$$

where C is a constant depending only on α_1 .

Finally, we apply Theorem 2. In the particular case that the X_{2j} , $j = 1, \dots, q-1$, are independent, one can show that (see (D.1))

$$r(x_1, x_2) = \sum_{j=1}^{q-1} \frac{p_{X_1, X_{2j}}(x_1, x_{2j})}{p_{X_1}(x_1)p_{X_{2j}}(x_{2j})} - (q-2),$$

where $p_{X_1, X_{2j}}$ denotes the joint density of (X_1, X_{2j}) . In this particular case, condition (3.7) is thus much weaker than condition (3.6): we only need smoothness conditions on several kernels which involve only one- and two-dimensional design densities. In the general case, we have a similar result. Suppose that for each fixed x_1

$$\frac{p_{X_1, X_{2k}}(x_1, x_{2k})}{p_{X_1}(x_1)p_{X_{2k}}(x_{2k})} \in \mathcal{H}(\beta, h_{1k}(x_1)) \quad (4.17)$$

with $h_{1k} \in L^2(\mathbb{P}^{X_1})$, for all $k = 1, \dots, q-1$. Moreover, suppose that for each fixed x_{2j}

$$\frac{p_{X_{2j}, X_{2k}}(x_{2j}, x_{2k})}{p_{X_{2j}}(x_{2j})p_{X_{2k}}(x_{2k})} \in \mathcal{H}(\beta, h'_{jk}(x_{2j})) \quad (4.18)$$

with $h'_{jk} \in L^2(\mathbb{P}^{X_{2j}})$, for all $j, k = 1, \dots, q-1$. Then we have:

Corollary 7. *Let $q > 2$. Let Assumption 2, 9, 10, and 4 be satisfied. Suppose that (4.17) and (4.18) are satisfied. Moreover, suppose that $\alpha_2 + \beta/(2\alpha_1 + 1) > \alpha_1/(2\alpha_1 + 1)$. Then (4.16) holds, where C is a constant depending only on α_1 .*

A proof of Corollary 7 is given in Appendix D. Note that in Corollary 7, the smoothness condition is stronger than the smoothness condition given in (4.15).

5. PROOF OF THEOREM 1 AND 2

5.1. The finite sample geometry. In this section, we present an empirical version of Assumption 2, which holds on the event

$$\mathcal{E}_\delta = \{ (1 - \delta)\|g\|^2 \leq \|g\|_n^2 \leq (1 + \delta)\|g\|^2 \text{ for all } g \in V \},$$

where $0 < \delta < 1$ is the constant from Theorem 1. Moreover, using concentration of measure inequalities for structured random matrices, we lower bound the probability that the event \mathcal{E}_δ occurs.

In order to state our first result, we introduce the empirical inner product $\langle \cdot, \cdot \rangle_n$ and the corresponding empirical norm $\| \cdot \|_n$ which are given by

$$\langle g, h \rangle_n = \frac{1}{n} \sum_{i=1}^n g(X^i)h(X^i)$$

and $\|g\|_n^2 = \langle g, g \rangle_n$ for $g, h \in L^2(\mathbb{P}^X)$.

Proposition 1. *Let Assumptions 1 and 2 hold. If \mathcal{E}_δ holds, then we have*

$$\|g_1 + g_2\|_n^2 \geq \frac{(1-\delta)}{(1+\delta)}(1-\rho_0)(\|g_1\|_n^2 + \|g_2\|_n^2) \quad \text{and} \quad (5.1)$$

$$\|g_1 + g_2\|_n^2 \geq \frac{(1-\delta)}{(1+\delta)}(1-\rho_0^2)\|g_1\|_n^2 \quad (5.2)$$

for all $g_1 \in V_1$, $g_2 \in V_2$, and also

$$\frac{|\langle g_1, g_2 \rangle_n|}{\|g_1\|_n \|g_2\|_n} \leq 1 - \frac{(1-\delta)}{(1+\delta)}(1-\rho_0) \quad (5.3)$$

for all $0 \neq g_1 \in V_1$, $0 \neq g_2 \in V_2$.

Proof. Let $g_1 \in V_1$ and $g_2 \in V_2$. By the definition of \mathcal{E}_δ and Lemma 1 combined with Assumption 2, we have

$$\begin{aligned} \|g_1 + g_2\|_n^2 &\geq (1-\delta)\|g_1 + g_2\|^2 \geq (1-\delta)(1-\rho_0)(\|g_1\|^2 + \|g_2\|^2) \\ &\geq \frac{(1-\delta)}{(1+\delta)}(1-\rho_0)(\|g_1\|_n^2 + \|g_2\|_n^2) \end{aligned}$$

and similarly

$$\begin{aligned} \|g_1 + g_2\|_n^2 &\geq (1-\delta)\|g_1 + g_2\|^2 \geq (1-\delta)(1-\rho_0^2)\|g_1\|^2 \\ &\geq \frac{(1-\delta)}{(1+\delta)}(1-\rho_0^2)\|g_1\|_n^2. \end{aligned}$$

This gives (5.1) and (5.2). (5.3) follows from (5.1) and Lemma 1. This completes the proof. \square

The following result follows from Rauhut [26, Theorem 7.3] (see also [29, Theorem 3.1]). It can also be obtained from a combination of Talagrand's inequality and Rudelson's lemma (see [28, Theorem 1]).

Theorem 7. *Let Assumption 3 hold. Then we have*

$$\mathbb{P}(\mathcal{E}_\delta) \geq 1 - 2^{3/4}d \exp\left(-\kappa \frac{n\delta^2}{\varphi^2 d}\right),$$

where κ is a universal constant.

Proof. Let b_1, \dots, b_d be an orthonormal basis of V . By Assumption 3 and [5, Lemma 1], we have

$$\left\| \sum_{j=1}^d b_j^2 \right\|_\infty \leq \varphi^2 d. \quad (5.4)$$

Now let

$$B_n = (\langle b_j, b_k \rangle_n)_{1 \leq j, k \leq d}.$$

Then [26, Theorem 7.3] (in the case $s = N = d$) yields for $0 < \delta < 1$,

$$\mathbb{P}(\|B_n - I\|_{\text{op}} \leq \delta) \geq 1 - 2^{3/4}d \exp\left(-\kappa \frac{n\delta^2}{\varphi^2 d}\right), \quad (5.5)$$

where $\kappa > 0$ is a universal constant. Here, $\|\cdot\|_{\text{op}}$ denotes the operator norm (see Lemma 6 (iii) for the definition). Note that we can apply [26, Theorem 7.3] since in the proof the condition [26, (4.2)] is only used in the form [26, (7.5)], which is satisfied by (5.4). A similar result follows from [29, Theorem 3.1].

Now, a function $g \in V$ with $\|g\| \leq 1$ can be written uniquely as $g = \sum_{j=1}^d x_j b_j$ with $x \in \mathbb{R}^d$ and $\|x\|_2 \leq 1$. Using this and $\|g\|_n^2 = x^T B_n x$, we obtain

$$\sup_{g \in V, \|g\| \leq 1} |\|g\|_n^2 - \|g\|^2| = \sup_{x \in \mathbb{R}^d, \|x\|_2 \leq 1} |x^T (B_n - I)x| = \|B_n - I\|_{\text{op}}, \quad (5.6)$$

where the latter equality follows from the spectral theorem. Moreover, we have that \mathcal{E}_δ holds if and only if

$$\sup_{g \in V, \|g\| \leq 1} |\|g\|_n^2 - \|g\|^2| \leq \delta.$$

Applying this, (5.5), and (5.6), we complete the proof. \square

5.2. Analysis of the variance via Von Neumann's theorem. The basic theorem in the theory of projections on sumspaces is due to von Neumann [34]. We state the following version dealing only with the first component, which is a consequence of [1, (15) on page 378] (see also [3, Theorem 2.C in Appendix A.4]).

Lemma 5. *Let \mathcal{H}_1 and \mathcal{H}_2 be two closed subspaces of a Hilbert space \mathcal{H} with inner product $\langle \cdot, \cdot \rangle$ and norm $\|\cdot\|$. Suppose that $\rho_0(\mathcal{H}_1, \mathcal{H}_2) < 1$. Let Π , Π_1 , Π_2 be the orthogonal projections on $\mathcal{H}_1 + \mathcal{H}_2$, \mathcal{H}_1 , \mathcal{H}_2 ,*

respectively. Let $h \in \mathcal{H}$, and let $(\Pi h)_1 \in \mathcal{H}_1$, $(\Pi h)_2 \in \mathcal{H}_2$ be the unique elements such that $\Pi h = (\Pi h)_1 + (\Pi h)_2$. Then

$$\left\| (\Pi h)_1 - \left(\Pi_1 - \sum_{j=1}^k (\Pi_1 \Pi_2)^j (1 - \Pi_1) \right) h \right\| \rightarrow 0$$

as k goes to infinity.

Remark 7. If we set $h_1^{(1)} = \Pi_1 h$ and proceed iteratively by setting $h_2^{(m)} = \Pi_2(h - h_1^{(m)})$ and $h_1^{(m+1)} = \Pi_1(h - h_2^{(m)})$, $m \geq 1$, then Lemma 5 can be rewritten as $\|(\Pi h)_1 - h_1^{(m)}\| \rightarrow 0$. This procedure is often called “backfitting”.

In this section, we apply Lemma 5 to the finite sample setting, using results from the previous section. Recall that $\hat{\Pi}_V$ is the orthogonal projection from \mathbb{R}^n to the subspace $\{(g(X^1), \dots, g(X^n))^T | g \in V\}$, and that $\hat{\Pi}_{V_1}$, $\hat{\Pi}_{V_2}$, and $\hat{\Pi}_{W_1}$ are defined analogously (replace V by V_1 , V_2 , and W_1 , respectively). We first prove:

Proposition 2. *Let Assumption 1 and 2 be satisfied. Let*

$$\rho_{0,\delta} = 1 - \frac{(1 - \delta)}{(1 + \delta)}(1 - \rho_0).$$

If \mathcal{E}_δ holds, then we have

$$\mathbb{E} \left[\|\hat{\Pi}_{W_1}(\hat{\Pi}_V \epsilon)_1\|_n^2 | X^1, \dots, X^n \right] \leq \frac{\sigma^2 \dim W_1}{n} + \frac{1}{1 - \rho_{0,\delta}^2} \frac{\sigma^2 \text{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2})}{n}.$$

In the proof, we will need the following result:

Lemma 6. *Let $A \in \mathbb{R}^{k_1 \times k_2}$ and $B \in \mathbb{R}^{k_2 \times k_1}$. Then*

- (i) $\text{tr}(AB) = \text{tr}(BA)$.
- (ii) $|\text{tr}(AB)| \leq \sqrt{\text{tr}(AA^T) \text{tr}(BB^T)}$.
- (iii) *Let $k_1 = k_2$ and B be symmetric and positive semi-definite. Then*

$$|\text{tr}(AB)| \leq \|A\|_{\text{op}} \text{tr}(B),$$

where $\|A\|_{\text{op}} = \sup_{\|x\|_2=1} \|Ax\|_2$ denotes the operator norm. Here, $\|\cdot\|_2$ denotes the Euclidean norm.

For completeness, a proof of Lemma 6 (iii) is given in Appendix E.

Proof of Proposition 2. Throughout the proof, suppose that \mathcal{E}_δ holds. Furthermore, we consider V as a subset of \mathbb{R}^n . This is no restriction, since (on \mathcal{E}_δ) each element $g \in V$ is uniquely determined by

$(g(X^1), \dots, g(X^n))^T$. From (5.3) and Lemma 5 applied to $(\mathcal{H}, \langle \cdot, \cdot \rangle) = (V, \langle \cdot, \cdot \rangle_n)$ and $\mathcal{H}_j = V_j$, $j = 1, 2$, we have

$$\|(\hat{\Pi}_V \epsilon)_1 - (\hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1})) \epsilon\|_n \rightarrow 0$$

as k goes to infinity. From (5.3), we have

$$\|\hat{\Pi}_{V_1} g_2\|_n \leq \rho_{0,\delta} \|g_2\|_n \quad (5.7)$$

for all $g_2 \in V_2$, which follows from

$$\|\hat{\Pi}_{V_1} g_2\|_n^2 = \langle \hat{\Pi}_{V_1} g_2, \hat{\Pi}_{V_1} g_2 \rangle_n = \langle \hat{\Pi}_{V_1} g_2, g_2 \rangle_n \leq \rho_{0,\delta} \|\hat{\Pi}_{V_1} g_2\|_n \|g_2\|_n.$$

Similarly, we have

$$\|\hat{\Pi}_{V_2} g_1\|_n \leq \rho_{0,\delta} \|g_1\|_n \quad (5.8)$$

for all $g_1 \in V_1$. This gives the improved convergence result

$$\begin{aligned} & \|(\hat{\Pi}_V \epsilon)_1 - (\hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1})) \epsilon\|_n \\ & \leq \sum_{j=k+1}^{\infty} \|(\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1}) \epsilon\|_n \\ & \leq \sum_{j=k+1}^{\infty} \rho_{0,\delta}^{2j-1} \|\epsilon\|_n \\ & = \frac{\rho_{0,\delta}^{2k+1}}{1 - \rho_{0,\delta}^2} \|\epsilon\|_n, \end{aligned}$$

and also

$$\begin{aligned} & \|\hat{\Pi}_{W_1} (\hat{\Pi}_V \epsilon)_1 - \hat{\Pi}_{W_1} (\hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1})) \epsilon\|_n \\ & \leq \frac{\rho_{0,\delta}^{2k+1}}{1 - \rho_{0,\delta}^2} \|\epsilon\|_n. \end{aligned} \quad (5.9)$$

Applying (5.9) and the bound $(x+y)^2 \leq (1+\epsilon)x^2 + (1+1/\epsilon)y^2$, $\epsilon > 0$, and then taking expectation, we obtain

$$\begin{aligned} & \mathbb{E} \left[\|\hat{\Pi}_{W_1}(\hat{\Pi}_V \epsilon)_1\|_n^2 \mid X^1, \dots, X^n \right] \\ & \leq (1+\epsilon) \mathbb{E} \left[\left\| \hat{\Pi}_{W_1} \left(\hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1}) \right) \epsilon \right\|_n^2 \mid X^1, \dots, X^n \right] \\ & \quad + (1+1/\epsilon) \frac{\rho_{0,\delta}^{4k+2}}{(1 - \rho_{0,\delta}^2)^2} \sigma^2. \end{aligned} \quad (5.10)$$

Since $\mathbb{E} [\|A\epsilon\|_n^2] = \sigma^2 \text{tr}(AA^T)/n$ for all $A \in \mathbb{R}^{n \times n}$, it remains to bound the trace of

$$\hat{\Pi}_{W_1} \left(\hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1}) \right) \left(\hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1}) \right)^T \hat{\Pi}_{W_1}. \quad (5.11)$$

Using

$$\begin{aligned} & \hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1}) \\ & = \sum_{j=0}^{k-1} \hat{\Pi}_{V_1} (\hat{\Pi}_{V_2} \hat{\Pi}_{V_1})^j (1 - \hat{\Pi}_{V_2}) + (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^k \hat{\Pi}_{V_1}, \end{aligned}$$

(5.11) is equal to

$$\begin{aligned} & \hat{\Pi}_{W_1} \sum_{j=0}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j \hat{\Pi}_{V_1} \left(\hat{\Pi}_{V_1} - (1 - \hat{\Pi}_{V_1}) \sum_{j=1}^k (\hat{\Pi}_{V_2} \hat{\Pi}_{V_1})^j \right) \hat{\Pi}_{W_1} \\ & - \hat{\Pi}_{W_1} \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j \left((1 - \hat{\Pi}_{V_2}) \sum_{j=0}^{k-1} (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j \hat{\Pi}_{V_1} + (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^k \hat{\Pi}_{V_1} \right) \hat{\Pi}_{W_1} \end{aligned}$$

and, since $\hat{\Pi}_{V_1}(1 - \hat{\Pi}_{V_1}) = 0$ and $\hat{\Pi}_{V_2}(1 - \hat{\Pi}_{V_2}) = 0$, this is equal to

$$\hat{\Pi}_{W_1} \sum_{j=0}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j \hat{\Pi}_{W_1} - \hat{\Pi}_{W_1} \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^{j+k} \hat{\Pi}_{W_1}.$$

By Lemma 6 (i) and the identities $\hat{\Pi}_{W_1} \hat{\Pi}_{V_1} = \hat{\Pi}_{W_1}$ and $\hat{\Pi}_{V_2} \hat{\Pi}_{V_2} = \hat{\Pi}_{V_2}$, we have for $j = 1, \dots, 2k$,

$$\begin{aligned} \text{tr}(\hat{\Pi}_{W_1} (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j \hat{\Pi}_{W_1}) & = \text{tr}((\hat{\Pi}_{V_2} \hat{\Pi}_{V_1})^{j-1} \hat{\Pi}_{V_2} \hat{\Pi}_{W_1} \hat{\Pi}_{V_2}) \\ & = \text{tr}((\hat{\Pi}_{V_2} \hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^{j-1} \hat{\Pi}_{V_2} \hat{\Pi}_{W_1} \hat{\Pi}_{V_2}). \end{aligned}$$

Thus the trace of (5.11) is bounded by

$$\dim W_1 + \sum_{j=1}^{2k} |\operatorname{tr}((\hat{\Pi}_{V_2} \hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^{j-1} \hat{\Pi}_{V_2} \hat{\Pi}_{W_1} \hat{\Pi}_{V_2})|. \quad (5.12)$$

Applying Lemma 6 (iii), this can be bounded by

$$\dim W_1 + \sum_{j=0}^{2k-1} \|\hat{\Pi}_{V_2} \hat{\Pi}_{V_1} \hat{\Pi}_{V_2}\|_{\text{op}}^j \operatorname{tr}(\hat{\Pi}_{V_2} \hat{\Pi}_{W_1} \hat{\Pi}_{V_2}).$$

By (5.7) and (5.8), we have

$$\|\hat{\Pi}_{V_2} \hat{\Pi}_{V_1} \hat{\Pi}_{V_2}\|_{\text{op}} \leq \rho_{0,\delta}^2. \quad (5.13)$$

Moreover, we have $\operatorname{tr}(\hat{\Pi}_{V_2} \hat{\Pi}_{W_1} \hat{\Pi}_{V_2}) = \operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2})$. Thus (5.12) is bounded by

$$\dim W_1 + \frac{1}{1 - \rho_{0,\delta}^2} \operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2}). \quad (5.14)$$

From (5.10)-(5.14), we conclude that

$$\begin{aligned} & \mathbb{E} \left[\|\hat{\Pi}_{W_1} (\hat{\Pi}_V \epsilon)_1\|_n^2 | X^1, \dots, X^n \right] \\ & \leq (1 + \epsilon) \left(\frac{\sigma^2 \dim W_1}{n} + \frac{1}{1 - \rho_{0,\delta}^2} \frac{\sigma^2 \operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2})}{n} \right) \\ & \quad + (1 + 1/\epsilon) \frac{\rho_{0,\delta}^{4k+2}}{(1 - \rho_{0,\delta}^2)^2} \sigma^2. \end{aligned}$$

Now, send k off to infinity first, and then let ϵ go to zero. This completes the proof. \square

From Proposition 2, we obtain a first upper bound for the variance term, which does not depend on the dimension of V_2 . Note that in Proposition 3 and Corollary 9, we show that this upper bound can be further refined.

Corollary 8. *Let Assumption 1 and 2 be satisfied. If \mathcal{E}_δ holds, then we have*

$$\mathbb{E} \left[\|\hat{\Pi}_{W_1} (\hat{\Pi}_V \epsilon)_1\|_n^2 | X^1, \dots, X^n \right] \leq \frac{1 + \delta}{1 - \delta} \frac{1}{1 - \rho_0^2} \frac{\sigma^2 \dim W_1}{n}.$$

Proof. By Lemma 6 (i) and (ii), we have

$$\operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2}) = \operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2} \hat{\Pi}_{W_1} \hat{\Pi}_{W_1}).$$

Applying Lemma 6 (iii) and (5.13), we obtain on \mathcal{E}_δ ,

$$\operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2}) \leq \|\hat{\Pi}_{W_1} \hat{\Pi}_{V_2} \hat{\Pi}_{W_1}\|_{\text{op}} \operatorname{tr}(\hat{\Pi}_{W_1}) \leq \rho_{0,\delta}^2 \dim W_1.$$

Thus, if \mathcal{E}_δ holds, Proposition 2 yields

$$\mathbb{E} \left[\|\hat{\Pi}_{W_1}(\hat{\Pi}_V \epsilon)_1\|_n^2 | X^1, \dots, X^n \right] \leq \frac{1}{1 - \rho_{0,\delta}^2} \frac{\sigma^2 \dim W_1}{n}.$$

Since $(1 - c(1 - \varrho))^2 \leq 1 - c(1 - \varrho^2)$ for $\varrho \in [0, 1]$ and a constant $0 \leq c \leq 1$ (both functions are equal to 1 at the right endpoint $\varrho = 1$ and the derivative of the left hand side is greater or equal than the derivative of the right hand side for all $\varrho \in [0, 1]$), we obtain (set $c = (1 - \delta)/(1 + \delta)$ and $\varrho = \rho_0$)

$$\frac{1}{1 - \rho_{0,\delta}^2} \leq \frac{1 + \delta}{1 - \delta} \frac{1}{1 - \rho_0^2}. \quad (5.15)$$

This completes the proof. \square

Proposition 3. *Let Assumption 1, 2, and 3 be satisfied. Then*

$$\frac{1}{n} \mathbb{E} \left[1_{\mathcal{E}_\delta} \operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2}) \right] \leq \frac{1}{(1 - \delta)^2} \left(\frac{\|\Pi_{V_2}|_{W_1}\|_{HS}^2}{n} + \frac{\dim W_1}{n} \frac{\varphi^2 d}{n} \right).$$

Remark 8. By considering $\Pi_{W_1} \Pi_{V_2} \Pi_{W_1}$ as a map from $W_1 \subseteq L^2(\mathbb{P}^{X_1})$ to itself, we have $\|\Pi_{V_2}|_{W_1}\|_{HS}^2 = \operatorname{tr}(\Pi_{W_1} \Pi_{V_2} \Pi_{W_1})$.

Proof. Let $\{\phi_{1j}\}_{1 \leq j \leq \dim W_1}$ be an orthonormal basis of W_1 , and let $\{\phi_{2j}\}_{1 \leq j \leq d_2}$ be an orthonormal basis of V_2 . Let

$$Z_1 = (\phi_{1j}(X_1^i))_{1 \leq i \leq n, 1 \leq j \leq \dim W_1}$$

and

$$Z_2 = (\phi_{2j}(X_2^i))_{1 \leq i \leq n, 1 \leq j \leq d_2},$$

Now, suppose that \mathcal{E}_δ holds. Then we have $\hat{\Pi}_{W_1} = Z_1(Z_1^T Z_1)^{-1} Z_1^T$ and $\hat{\Pi}_{V_2} = Z_2(Z_2^T Z_2)^{-1} Z_2^T$. Thus

$$\operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2}) = \operatorname{tr} \left(\left(\frac{1}{n} Z_1^T Z_1 \right)^{-1} \frac{1}{n} Z_1^T Z_2 \left(\frac{1}{n} Z_2^T Z_2 \right)^{-1} \frac{1}{n} Z_2^T Z_1 \right),$$

where we applied Lemma 6 (i). By Theorem 7, we have $\|(1/n) Z_j^T Z_j - I\|_{\text{op}} \leq \delta$ and thus

$$\left(\frac{1}{n} Z_j^T Z_j \right)^{-1} = \sum_{k \geq 0} \left(I - \frac{1}{n} Z_j^T Z_j \right)^k$$

for $j = 1, 2$. We conclude that

$$\begin{aligned}
& \mathbb{E} \left[1_{\mathcal{E}_\delta} \operatorname{tr}(\hat{\Pi}_{W_1} \hat{\Pi}_{V_2}) \right] \\
& \leq \sum_{k,l=0}^{\infty} \mathbb{E} \left[1_{\mathcal{E}_\delta} \left| \operatorname{tr} \left(\left(I - \frac{1}{n} Z_1^T Z_1 \right)^k \frac{1}{n} Z_1^T Z_2 \left(I - \frac{1}{n} Z_2^T Z_2 \right)^l \frac{1}{n} Z_2^T Z_1 \right) \right| \right] \\
& \leq \sum_{k,l=0}^{\infty} \mathbb{E} \left[1_{\mathcal{E}_\delta} \sqrt{\operatorname{tr} \left(\left(I - \frac{1}{n} Z_1^T Z_1 \right)^{2k} \frac{1}{n} Z_1^T Z_2 \left(\frac{1}{n} Z_1^T Z_2 \right)^T \right)} \right. \\
& \quad \cdot \left. \sqrt{\operatorname{tr} \left(\left(I - \frac{1}{n} Z_2^T Z_2 \right)^{2l} \frac{1}{n} Z_2^T Z_1 \left(\frac{1}{n} Z_2^T Z_1 \right)^T \right)} \right] \\
& \leq \sum_{k,l=0}^{\infty} \delta^{k+l} \mathbb{E} \left[\operatorname{tr} \left(\frac{1}{n} Z_1^T Z_2 \left(\frac{1}{n} Z_1^T Z_2 \right)^T \right) \right],
\end{aligned}$$

where we applied Lemma 6 (i) and (ii) in the second inequality and Lemma 6 (i) and (iii) and the bound $\|(1/n)Z_j^T Z_j - I\|_{\text{op}} \leq \delta$ in the third inequality. Now

$$\begin{aligned}
& \mathbb{E} \left[\operatorname{tr} \left(\frac{1}{n} Z_1^T Z_2 \left(\frac{1}{n} Z_1^T Z_2 \right)^T \right) \right] \\
& = \sum_{j=1}^{\dim W_1} \sum_{k=1}^{d_2} \left(\mathbb{E} [\phi_{1j}(X_1) \phi_{2k}(X_2)]^2 + \frac{1}{n} \operatorname{Var}(\phi_{1j}(X_1) \phi_{2k}(X_2)) \right) \\
& \leq \sum_{j=1}^{\dim W_1} \sum_{k=1}^{d_2} \langle \phi_{1j}, \phi_{2k} \rangle^2 + \dim W_1 \frac{\varphi^2 d}{n},
\end{aligned}$$

where we applied (5.4). Finally, we use

$$\|\Pi_{V_2}|_{W_1}\|_{HS}^2 = \sum_{j=1}^{\dim W_1} \|\Pi_{V_2} \phi_{1j}\|^2 = \sum_{j=1}^{\dim W_1} \sum_{k=1}^{d_2} \langle \phi_{1j}, \phi_{2k} \rangle^2$$

This completes the proof. \square

Combining Proposition 2 and 3, we obtain the following improvement of Corollary 8:

Corollary 9. *Let Assumption 1, 2, and 3 be satisfied. Then*

$$\begin{aligned} & \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\hat{\Pi}_{W_1}(\hat{\Pi}_V \epsilon)_1\|_n^2 \right] \\ & \leq \frac{\sigma^2 \dim W_1}{n} + \frac{(1+\delta)}{(1-\delta)^3} \frac{1}{1-\rho_0^2} \left(\frac{\sigma^2 \|\Pi_{V_2}|_{W_1}\|_{HS}^2}{n} + \frac{\sigma^2 \dim W_1}{n} \frac{\varphi^2 d}{n} \right). \end{aligned}$$

5.3. End of the proof of Theorem 1. Applying the arguments of [2], we obtain

$$E \left[\|f_1 - \hat{f}_1^*\|^2 \right] \leq \mathbb{E} \left[1_{\mathcal{E}_\delta} \|f_1 - (\hat{f}_V)_1\|^2 \right] + R_n. \quad (5.16)$$

The details can be found in Appendix F. By the projection theorem (see, e.g., [27, Theorem II.3]), we have

$$\|f_1 - (\hat{f}_V)_1\|^2 = \|f_1 - \Pi_{V_1} f_1\|^2 + \|\Pi_{V_1} f_1 - (\hat{f}_V)_1\|^2. \quad (5.17)$$

By the definition of \mathcal{E}_δ and the definition of \hat{f}_V , we have

$$\begin{aligned} & \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\Pi_{V_1} f_1 - (\hat{f}_V)_1\|_n^2 \right] \\ & \leq \frac{1}{(1-\delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\Pi_{V_1} f_1 - (\hat{f}_V)_1\|_n^2 \right] \\ & = \frac{1}{(1-\delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\Pi_{V_1} f_1 - (\hat{\Pi}_V \mathbf{Y})_1\|_n^2 \right] \\ & = \frac{1}{(1-\delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \mathbb{E} \left[\|\Pi_{V_1} f_1 - (\hat{\Pi}_V \mathbf{Y})_1\|_n^2 | X^1, \dots, X^n \right] \right]. \quad (5.18) \end{aligned}$$

Recall from Section 2.3 that on \mathcal{E}_δ each $g \in V$ is determined uniquely by $(g(X^1), \dots, g(X^n))^T$, which implies that on \mathcal{E}_δ we don't have to distinguish between those objects. For instance, if \mathcal{E}_δ holds, then $\hat{\Pi}_V$ and $\hat{\Pi}_{V_1}$ can also be seen as maps from $L^2(\mathbb{P}^X)$ to V and V_1 , respectively (by letting $\hat{\Pi}_V h$ (resp. $\hat{\Pi}_{V_1} h$) equal to $\hat{\Pi}_V(h(X^1), \dots, h(X^n))^T$ (resp. $\hat{\Pi}_{V_1}(h(X^1), \dots, h(X^n))^T$) for $h \in L^2(\mathbb{P}^X)$). Moreover, if \mathcal{E}_δ holds, then we have $\mathbb{E}[(\hat{\Pi}_V \epsilon)_1 | X^1, \dots, X^n] = 0$ and $(\hat{\Pi}_V \mathbf{Y})_1 = (\hat{\Pi}_V f)_1 + (\hat{\Pi}_V \epsilon)_1$ (the former follows for instance from (5.9)). Thus (5.18) is equal to

$$\frac{1}{(1-\delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \left(\|\Pi_{V_1} f_1 - (\hat{\Pi}_V f)_1\|_n^2 + \mathbb{E} \left[\|(\hat{\Pi}_V \epsilon)_1\|_n^2 | X_1, \dots, X_n \right] \right) \right]. \quad (5.19)$$

From (5.17)-(5.19) and the definition of \mathcal{E}_δ , we obtain

$$\begin{aligned} & \mathbb{E} \left[1_{\mathcal{E}_\delta} \|f_1 - (\hat{f}_V)_1\|^2 \right] \\ & \leq \frac{(1+\delta)}{(1-\delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \left(\|f_1 - \Pi_{V_1} f_1\|^2 + \|\Pi_{V_1} f_1 - (\hat{\Pi}_V f)_1\|^2 \right) \right] \\ & \quad + \frac{1}{(1-\delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \mathbb{E} \left[\|(\hat{\Pi}_V \epsilon)_1\|_n^2 | X^1, \dots, X^n \right] \right]. \end{aligned} \quad (5.20)$$

Applying the projection theorem and Corollary 8, this is bounded by

$$\frac{(1+\delta)}{(1-\delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|f_1 - (\hat{\Pi}_V f)_1\|^2 \right] + \frac{(1+\delta)}{(1-\delta)^2} \frac{1}{1-\rho_0^2} \frac{\sigma^2 \dim V_1}{n}.$$

By Assumption 2 and Lemma 1, we have

$$\|f_1 - (\hat{\Pi}_V f)_1\|^2 \leq \frac{1}{1-\rho_0^2} \|f - \hat{\Pi}_V f\|^2.$$

The projection theorem implies that

$$\|f - \hat{\Pi}_V f\|^2 = \|f - \Pi_V f\|^2 + \|\Pi_V f - \hat{\Pi}_V f\|^2.$$

Now, the following proposition completes the proof.

Lemma 7. *Let Assumption 3 be satisfied. Then*

$$\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\Pi_V - \hat{\Pi}_V) f\|^2 \right] \leq \frac{1}{(1-\delta)^2} \frac{\varphi^2 d}{n} \|f - \Pi_V f\|^2.$$

Remark 9. Instead of applying Lemma 7, one can also apply the easier and weaker bound

$$\|\Pi_V f - \hat{\Pi}_V f\|^2 \leq \frac{1}{1-\delta} \|\Pi_V f - \hat{\Pi}_V f\|_n^2 \leq \frac{1}{1-\delta} \|f - \Pi_V f\|_n^2,$$

which follows from the definition of \mathcal{E}_δ and the projection theorem.

Proof of Lemma 7. Throughout the proof suppose that the event \mathcal{E}_δ holds. Let b_1, \dots, b_d be an orthonormal basis of V . Let

$$B_n = (\langle b_j, b_k \rangle_n)_{1 \leq j, k \leq d},$$

and let

$$x = (\langle b_1, f \rangle, \dots, \langle b_d, f \rangle)^T \quad \text{and} \quad x_n = (\langle b_1, f \rangle_n, \dots, \langle b_d, f \rangle_n)^T.$$

Then we have

$$\Pi_V f = \sum_{j=1}^d x_j b_j \quad \text{and} \quad \hat{\Pi}_V f = \sum_{j=1}^d (B_n^{-1} x_n)_j b_j$$

and thus

$$\|(\Pi_V - \hat{\Pi}_V) f\|^2 = \|B_n^{-1} x_n - x\|_2^2. \quad (5.21)$$

Since \mathcal{E}_δ holds, we have $\|B_n - I\|_{\text{op}} \leq \delta$ (see the proof of Theorem 7). This implies that

$$B_n^{-1} = \sum_{k \geq 0} (I - B_n)^k.$$

Thus

$$\begin{aligned} B_n^{-1}x_n - x &= \sum_{k \geq 0} (I - B_n)^k x_n - x \\ &= \sum_{k \geq 0} (I - B_n)^k (x_n - x) + \sum_{k \geq 1} (I - B_n)^k x. \end{aligned}$$

Applying the bounds $\|B_n - I\|_{\text{op}} \leq \delta$ and $(x + y)^2 \leq (1 + \epsilon)x^2 + (1 + 1/\epsilon)y^2$, $\epsilon > 0$ arbitrary, we obtain

$$\begin{aligned} \|B_n^{-1}x_n - x\|_2^2 &\leq \frac{1}{(1 - \delta)^2} (\|x_n - x\|_2 + \|(B_n - I)x\|_2)^2 \\ &\leq \frac{1}{(1 - \delta)^2} ((1 + \epsilon)\|x_n - x\|_2^2 + (1 + 1/\epsilon)\|(B_n - I)x\|_2^2). \end{aligned} \quad (5.22)$$

Now we have

$$\begin{aligned} \mathbb{E} [\|x_n - x\|_2^2] &= \mathbb{E} \left[\sum_{j=1}^d (\langle b_j, f \rangle_n - \langle b_j, f \rangle)^2 \right] \\ &= \frac{1}{n} \sum_{j=1}^d \text{Var}(b_j(X)f(X)) \leq \frac{\varphi^2 d}{n} \|f\|^2, \end{aligned} \quad (5.23)$$

where we applied (5.4). The j th coordinate of $(B_n - I)x$ is equal to

$$\langle b_j, \sum_{k=1}^d b_k x_k \rangle_n - \langle b_j, f \rangle = \langle b_j, \Pi_V f \rangle_n - \langle b_j, \Pi_V f \rangle \quad (5.24)$$

and thus

$$\begin{aligned} \mathbb{E} [\|(B_n - I)x\|_2^2] &= \mathbb{E} \left[\sum_{j=1}^d (\langle b_j, \Pi_V f \rangle_n - \langle b_j, \Pi_V f \rangle)^2 \right] \\ &= \frac{1}{n} \sum_{j=1}^{d_1} \text{Var}(b_j(X)\Pi_V f(X)) \leq \frac{\varphi^2 d}{n} \|\Pi_V f\|^2. \end{aligned} \quad (5.25)$$

Applying (5.21)-(5.25), we conclude that

$$\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\Pi_V - \hat{\Pi}_V)f\|^2 \right] \leq \frac{1}{(1 - \delta)^2} \frac{\varphi^2 d}{n} ((1 + \epsilon)\|f\|^2 + (1 + 1/\epsilon)\|\Pi_V f\|^2). \quad (5.26)$$

Finally, since Π_V and $\hat{\Pi}_V$ fix elements in V , we obtain $(\Pi_V - \hat{\Pi}_V)\Pi_V f = 0$ and thus $(\Pi_V - \hat{\Pi}_V)f = (\Pi_V - \hat{\Pi}_V)(1 - \Pi_V)f$. Combining this with (5.26), we see that

$$\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\Pi_V - \hat{\Pi}_V)f\|^2 \right] \leq \frac{1}{(1-\delta)^2} \frac{\varphi^2 d}{n} (1+\epsilon) \|f - \Pi_V f\|^2.$$

Since $\epsilon > 0$ is arbitrary, this completes the proof. \square

5.4. End of the proof of Theorem 2. Compared with the previous section, we modify our analysis of the bias term, which is based the following two lemmas. Moreover, we replace Corollary 9 by Corollary 8.

Lemma 8. *Let Assumption 1 and 2 be satisfied. If \mathcal{E}_δ holds, then we have*

$$\|\hat{\Pi}_{V_1} h_1 - (\hat{\Pi}_V h_1)_1\|_n^2 \leq \frac{(1+\delta)}{(1-\delta)} \frac{1}{(1-\rho_0^2)} \|h_1 - \hat{\Pi}_{V_1} h_1\|_n^2$$

for all $h_1 \in H_1$ and

$$\|(\hat{\Pi}_V h_2)_1\|_n^2 \leq \frac{(1+\delta)}{(1-\delta)} \frac{1}{(1-\rho_0^2)} \|h_2 - \hat{\Pi}_{V_2} h_2\|_n^2$$

for all $h_2 \in H_2$.

Proof of Lemma 8. By (5.2) and the fact that projections lower the norm, we have

$$\begin{aligned} \|\hat{\Pi}_{V_1} h_1 - (\hat{\Pi}_V h_1)_1\|_n^2 &= \|(\hat{\Pi}_V(h_1 - \hat{\Pi}_{V_1} h_1))_1\|_n^2 \\ &\leq \frac{(1+\delta)}{(1-\delta)} \frac{1}{(1-\rho_0^2)} \|h_1 - \hat{\Pi}_{V_1} h_1\|_n^2. \end{aligned}$$

This gives the first inequality. Since $\hat{\Pi}_V \hat{\Pi}_{V_2} h_2 = \hat{\Pi}_{V_2} h_2$ and $(\hat{\Pi}_{V_2} h_2)_1 = 0$, we have $(\hat{\Pi}_V h_2)_1 = (\hat{\Pi}_V(h_2 - \hat{\Pi}_{V_2} h_2))_1$. Using the previous arguments, we conclude that

$$\|(\hat{\Pi}_V h_2)_1\|_n^2 = \|(\hat{\Pi}_V(h_2 - \hat{\Pi}_{V_2} h_2))_1\|_n^2 \leq \frac{(1+\delta)}{(1-\delta)} \frac{1}{(1-\rho_0^2)} \|h_2 - \hat{\Pi}_{V_2} h_2\|_n^2.$$

This completes the proof. \square

Lemma 9. *Let Assumption 1 and 2 be satisfied. Then we have*

$$\begin{aligned} &\mathbb{E} \left[1_{\mathcal{E}_\delta} \|f_1 - (\hat{\Pi}_V f)_1\|_n^2 \right] \\ &\leq \frac{(1+\delta)}{(1-\delta)} \frac{3}{(1-\rho_0^2)} (\|f_1 - \Pi_{V_1} f_1\|^2 + \|f_2 - \Pi_{V_2} f_2\|^2). \end{aligned}$$

Proof of Lemma 9. By the projection theorem, we have

$$\|f_1 - (\hat{\Pi}_V f)_1\|_n^2 = \|f_1 - \hat{\Pi}_{V_1} f_1\|_n^2 + \|\hat{\Pi}_{V_1} f_1 - (\hat{\Pi}_V f_1)_1 - (\hat{\Pi}_V f_2)_1\|_n^2.$$

Applying this, the bound $\|x + y\|^2 \leq 2\|x\|^2 + 2\|y\|^2$, and Lemma 9, we obtain

$$\begin{aligned} & \mathbb{E} \left[1_{\mathcal{E}_\delta} \|f_1 - (\hat{\Pi}_V f)_1\|_n^2 \right] \\ & \leq \frac{(1 + \delta)}{(1 - \delta)} \frac{3}{(1 - \rho_0^2)} \left(\mathbb{E} \left[\|f_1 - \hat{\Pi}_{V_1} f_1\|_n^2 + \|f_2 - \hat{\Pi}_{V_2} f_2\|_n^2 \right] \right). \end{aligned}$$

By the projection theorem, we have for $j = 1, 2$,

$$\|f_j - \hat{\Pi}_{V_j} f_j\|_n^2 \leq \|f_j - \Pi_{V_j} f_j\|_n^2.$$

Moreover, taking expectation, we get for $j = 1, 2$,

$$\mathbb{E} [\|f_j - \Pi_{V_j} f_j\|_n^2] = \|f_j - \Pi_{V_j} f_j\|^2.$$

This completes the proof. \square

Now, we begin with the proof of Theorem 2. Repeating the steps (5.16)-(5.19) in the proof of Theorem 1, we have

$$\begin{aligned} & \mathbb{E} [\|f_1 - \hat{f}_1^*\|^2] \\ & \leq \|f_1 - \Pi_{W_1} f_1\|^2 + \frac{1}{(1 - \delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\Pi_{W_1} f_1 - \hat{\Pi}_{W_1} (\hat{\Pi}_V f)_1\|_n^2 \right] \\ & \quad + \frac{1}{(1 - \delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \mathbb{E} \left[\|\hat{\Pi}_{W_1} (\hat{\Pi}_V \epsilon)_1\|_n^2 | X^1, \dots, X^n \right] \right] + R_n. \end{aligned}$$

Applying the bound $\|x + y\|^2 \leq 2\|x\|^2 + 2\|y\|^2$ and the fact that projections lower the norm, we obtain

$$\begin{aligned} & \|\Pi_{W_1} f_1 - \hat{\Pi}_{W_1} (\hat{\Pi}_V f)_1\|_n^2 \\ & \leq 2\|\Pi_{W_1} f_1 - \hat{\Pi}_{W_1} f_1\|_n^2 + 2\|\hat{\Pi}_{W_1} f_1 - \hat{\Pi}_{W_1} (\hat{\Pi}_V f)_1\|_n^2 \\ & \leq 2\|\Pi_{W_1} f_1 - \hat{\Pi}_{W_1} f_1\|_n^2 + 2\|f_1 - (\hat{\Pi}_V f)_1\|_n^2. \end{aligned}$$

Thus

$$\begin{aligned} & \mathbb{E} [\|f_1 - \hat{f}_1^*\|^2] \\ & \leq \|f_1 - \Pi_{W_1} f_1\|^2 + \frac{2}{(1 - \delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\Pi_{W_1} f_1 - \hat{\Pi}_{W_1} f_1\|_n^2 \right] \\ & \quad + \frac{1}{(1 - \delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \mathbb{E} \left[\|\hat{\Pi}_{W_1} (\hat{\Pi}_V \epsilon)_1\|_n^2 | X^1, \dots, X^n \right] \right] \\ & \quad + \frac{2}{(1 - \delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|f_1 - (\hat{\Pi}_V f)_1\|_n^2 \right] + R_n. \end{aligned} \tag{5.27}$$

Similarly as in Lemma 7, we have

$$\mathbb{E} \left[1_{\mathcal{E}_\delta} \|\Pi_{W_1} f_1 - \hat{\Pi}_{W_1} f_1\|^2 \right] \leq \frac{1}{(1-\delta)^2} \frac{\varphi^2 d}{n} \|f_1 - \Pi_{W_1} f_1\|^2. \quad (5.28)$$

Inserting (5.28), Lemma 9, and Corollary 9 into (5.27), we complete the proof. \square

6. PROOF OF THEOREM 3 AND 4

6.1. The bias term revisited. In this subsection, we prove two results which lead to improvements of the bias term considered in Lemma 9. This improvement is possible under the additional regularity conditions on the design densities, namely (3.6) or (3.7). We first prove:

Proposition 4. *Let Assumption 1, 2, and 4 be satisfied. Let r , $\phi(V_2)$, $\psi_\Pi(V_2)$, and h_1 be as in Theorem 4. Then*

$$\|(\Pi_V f_2)_1\| \leq \frac{1}{1-\rho_0^2} \|h_1\| \psi_\Pi(V_2) \phi(V_2). \quad (6.1)$$

Proof. Let $\phi_1, \dots, \phi_{d_1}$ be an orthonormal basis of V_1 . We have

$$\begin{aligned} & \|\Pi_{V_1}(f_2 - \Pi_{V_2} f_2)\|^2 \\ &= \sum_{j=1}^{d_1} \left(\int_{S_1 \times S_2} \phi_j(x_1) (f_2(x_2) - \Pi_{V_2} f_2(x_2)) p(x_1, x_2) d(\nu_1 \otimes \nu_2)(x_1, x_2) \right)^2 \\ &= \sum_{j=1}^{d_1} \left(\int_{S_1} \phi_j(x_1) \int_{S_2} ((1 - \Pi_{V_2}) f_2(x_2)) \frac{p(x_1, x_2)}{p_1(x_1) p_2(x_2)} d\mathbb{P}^{X_2}(x_2) d\mathbb{P}^{X_1}(x_1) \right)^2 \\ &= \sum_{j=1}^{d_1} \left(\int_{S_1} \left\langle (1 - \Pi_{V_2}) f_2, \frac{p(x_1, \cdot)}{p_1(x_1) p_2(\cdot)} \right\rangle_{L^2(\mathbb{P}^{X_2})} \phi_j(x_1) d\mathbb{P}^{X_1}(x_1) \right)^2, \end{aligned}$$

where we already applied Fubini's theorem and Assumption 5 in the second equality. Since orthogonal projections are idempotent and self-adjoint, the above is equal to

$$\begin{aligned} &= \sum_{j=1}^{d_1} \left(\int_{S_1} \langle (1 - \Pi_{V_2}) f_2, (1 - \Pi_{V_2})(r(x_1, \cdot)) \rangle_{L^2(\mathbb{P}^{X_2})} \phi_j(x_1) d\mathbb{P}^{X_1}(x_1) \right)^2 \\ &\leq \int_{S_1} \langle (1 - \Pi_{V_2}) f_2, (1 - \Pi_{V_2})(r(x_1, \cdot)) \rangle_{L^2(\mathbb{P}^{X_2})}^2 d\mathbb{P}^{X_1}(x_1) \\ &\leq \|h_1\|^2 (\psi_\Pi(V_2) \phi(V_2))^2, \end{aligned}$$

where we applied Bessel's inequality in the first inequality and the Cauchy-Schwarz inequality and (3.7) in the second inequality. Thus

we have shown that

$$\|\Pi_{V_1}(f_2 - \Pi_{V_2}f_2)\| \leq \|h_1\|\psi_\Pi(V_2)\phi(V_2). \quad (6.2)$$

Now, by Lemma 5, we have

$$\|(\Pi_V h)_1 - (\Pi_{V_1} - \sum_{j=1}^k (\Pi_{V_1}\Pi_{V_2})^j(1 - \Pi_{V_1}))h\| \rightarrow 0, \quad (6.3)$$

as $k \rightarrow \infty$, for all $h \in L^2(\mathbb{P}^X)$. By Assumption 2, we have

$$\|\Pi_{V_1}\Pi_{V_2}h\| \leq \rho_0\|\Pi_{V_2}h\| \quad \text{and} \quad \|\Pi_{V_2}\Pi_{V_1}h\| \leq \rho_0\|\Pi_{V_1}h\|$$

for all $h \in L^2(\mathbb{P}^X)$, which follows as in the proof of (5.7). Applying this and (6.2), we obtain

$$\begin{aligned} & \|(\Pi_{V_1} - \sum_{j=1}^k (\Pi_{V_1}\Pi_{V_2})^j(1 - \Pi_{V_1}))(f_2 - \Pi_{V_2}f_2)\| \\ &= \left\| \sum_{j=0}^k (\Pi_{V_1}\Pi_{V_2})^j \Pi_{V_1}(f_2 - \Pi_{V_2}f_2) \right\| \\ &\leq \sum_{j=0}^k \rho_0^{2j} \|\Pi_{V_1}(f_2 - \Pi_{V_2}f_2)\| \leq \frac{1}{1 - \rho_0^2} \|h_1\|\psi_\Pi(V_2)\phi(V_2). \end{aligned} \quad (6.4)$$

Since $\Pi_V \Pi_{V_2}f_2 = \Pi_{V_2}f_2$ and $(\Pi_{V_2}f_2)_1 = 0$, we have $(\Pi_V f_2)_1 = (\Pi_V(f_2 - \Pi_{V_2}f_2))_1$. Applying this, (6.3), and (6.4), we conclude that

$$\|(\Pi_V f_2)_1\| \leq \frac{1}{1 - \rho_0^2} \|h_1\|\psi_\Pi(V_2)\phi(V_2).$$

This completes the proof. \square

Proposition 5. *Let Assumption 1, 2, 3, and 4 be satisfied. Let $\phi(V_2)$, $\psi(V_2)$, and h_1 be as in Theorem 3. Moreover, suppose that $\|g_1\|_\infty \leq \varphi\sqrt{d_1}\|g_1\|$ for all $g_1 \in V_1$. Then*

$$\begin{aligned} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\hat{\Pi}_V f_2)_1\|_n^2 \right] &\leq \frac{(1 + \delta)^2}{(1 - \delta)^3} \frac{2}{(1 - \rho_0^2)^2} \left(\|h_1\|^2 (\psi(V_2)\phi(V_2))^2 \right. \\ &\quad \left. + \frac{1}{n} \|h_1\|^2 \|(1 - \Pi_{V_2})f_2\|_\infty^2 \psi^2(V_2) + \phi^2(V_2) \frac{\varphi^2 d_1}{n} \right). \end{aligned} \quad (6.5)$$

Proof. The proof is similar to the proof of Proposition 4. Throughout the proof, suppose that the event \mathcal{E}_δ holds. Then $\hat{\Pi}_V$ is a well-defined map from $L^2(\mathbb{P}^X)$ to V . Let $\phi_1, \dots, \phi_{d_1}$ be an orthonormal basis of V_1 .

By repeating the arguments at the beginning of the proof of Lemma 7, we obtain

$$\|\hat{\Pi}_{V_1}(f_2 - \hat{\Pi}_{V_2}f_2)\|_n^2 \leq \frac{1}{1-\delta} \sum_{j=1}^{d_1} \langle \phi_j, (1 - \hat{\Pi}_{V_2})f_2 \rangle_n^2.$$

Thus

$$\begin{aligned} & \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\hat{\Pi}_{V_1}(f_2 - \hat{\Pi}_{V_2}f_2)\|_n^2 \right] \\ & \leq \frac{2}{(1-\delta)} \mathbb{E} \left[\sum_{j=1}^{d_1} \langle \phi_{j, \Pi_2}, (1 - \hat{\Pi}_{V_2})f_2 \rangle_n^2 \right] \end{aligned} \quad (6.6)$$

$$+ \frac{2}{(1-\delta)} \mathbb{E} \left[\sum_{j=1}^{d_1} \langle \phi_j - \phi_{j, \Pi_2}, (1 - \hat{\Pi}_{V_2})f_2 \rangle_n^2 \right], \quad (6.7)$$

where

$$\phi_{j, \Pi_2}(x_2) = \int \phi_j(x_1) \frac{p(x_1, x_2)}{p_1(x_1)p_2(x_2)} p_1(x_1) d\nu_1(x_1)$$

is the conditional expectation of $\phi_j(X_1)$ given $X_2 = x_2$ (for \mathbb{P}^{X_2} -almost all x_2 , by Assumption 4). The expectation in (6.6) is equal to

$$\begin{aligned} & \mathbb{E} \left[\sum_{j=1}^{d_1} \left(\int \left\langle \frac{p(x_1, \cdot)}{p_1(x_1)p_2(\cdot)}, (1 - \hat{\Pi}_{V_2})f_2 \right\rangle_n \phi_j(x_1) p_1(x_1) d\nu_1(x_1) \right)^2 \right] \\ & \leq \int \mathbb{E} \left[\left\langle \frac{p(x_1, \cdot)}{p_1(x_1)p_2(\cdot)}, (1 - \hat{\Pi}_{V_2})f_2 \right\rangle_n^2 \right] p_1(x_1) d\nu_1(x_1), \end{aligned}$$

where we applied Bessel's inequality and Fubini's theorem in the last inequality. Applying the fact that orthogonal projections are idempotent and self-adjoint and then the Cauchy-Schwarz inequality, this is

$$\leq \int \mathbb{E} \left[\left\| (1 - \hat{\Pi}_{V_2}) \frac{p(x_1, \cdot)}{p_1(x_1)p_2(\cdot)} \right\|_n^2 \left\| (1 - \hat{\Pi}_{V_2})f_2 \right\|_n^2 \right] p_1(x_1) d\nu_1(x_1).$$

Applying the projection theorem and then (3.6), this is

$$\begin{aligned} & \leq \int \mathbb{E} \left[\left\| (1 - \Pi_{V_2}) \frac{p(x_1, \cdot)}{p_1(x_1)p_2(\cdot)} \right\|_n^2 \left\| (1 - \Pi_{V_2})f_2 \right\|_n^2 \right] p_1(x_1) d\nu_1(x_1) \\ & \leq \frac{n-1}{n} \|h_1\|^2 (\psi(V_2)\phi(V_2))^2 + \frac{1}{n} \|h_1\|^2 \|(1 - \Pi_{V_2})f_2\|_\infty^2 (\psi(V_2))^2. \end{aligned}$$

Now we turn to the expectation in (6.7). We have

$$\begin{aligned} \mathbb{E} \left[\sum_{j=1}^{d_1} \langle \phi_j - \phi_{j,\Pi_2}, (1 - \hat{\Pi}_{V_2})f_2 \rangle_n^2 \right] &\leq \frac{\varphi^2 d_1}{n} \mathbb{E} \left[\|(1 - \hat{\Pi}_{V_2})f_2\|_n^2 \right] \\ &\leq \frac{\varphi^2 d_1}{n} \|(1 - \Pi_{V_2})f_2\|^2. \end{aligned}$$

To prove the first inequality, first note that the $((1 - \hat{\Pi}_{V_2})f_2)(X_2^i)$ depend only on X_2^1, \dots, X_2^n and we have

$$\mathbb{E} \left[(\phi_j - \phi_{j,\Pi_2})(X^i) | X_2^1, \dots, X_2^n, X_1^{i'} \right] = \mathbb{E} \left[(\phi_j - \phi_{j,\Pi_2})(X^i) | X_2^i \right] = 0$$

for $i \neq i'$. This implies that the nondiagonal terms vanish. Next, apply the inequalities $\mathbb{E}[(\phi_j - \phi_{j,\Pi_2})^2(X) | X_2] \leq \mathbb{E}[\phi_j^2(X_1) | X_2]$ and

$$\left\| \sum_{j=1}^{d_1} \phi_j^2 \right\|_\infty \leq \varphi^2 d_1,$$

which follows from the bound $\|g_1\|_\infty \leq \varphi \sqrt{d_1} \|g_1\|$, for all $g_1 \in V_1$, and [5, Lemma 1]. Thus we have shown that

$$\begin{aligned} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\hat{\Pi}_{V_1}(f_2 - \hat{\Pi}_{V_2}f_2)\|_n^2 \right] &\leq \frac{2}{(1 - \delta)} (\|h_1\|^2 (\psi(V_2)\phi(V_2))^2 \\ &\quad + \frac{1}{n} \|h_1\|^2 \|(1 - \Pi_{V_2})f_2\|_\infty^2 (\psi(V_2))^2 + \phi^2(V_2) \frac{\varphi^2 d_1}{n}) \quad (6.8) \end{aligned}$$

The remaining arguments are as in the proof of Proposition 4. From (5.3) and Lemma 5 we have

$$\|(\hat{\Pi}_V h)_1 - (\hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1}))h\|_n \rightarrow 0 \quad (6.9)$$

as $k \rightarrow \infty$, for all $h \in L^2(\mathbb{P}^X)$. Applying (5.7) and (5.8) as in (6.4), we obtain

$$\begin{aligned} &\|(\hat{\Pi}_{V_1} - \sum_{j=1}^k (\hat{\Pi}_{V_1} \hat{\Pi}_{V_2})^j (1 - \hat{\Pi}_{V_1}))(f_2 - \hat{\Pi}_{V_2}f_2)\|_n \\ &\leq \frac{1}{1 - \rho_{0,\delta}^2} \|\hat{\Pi}_{V_1}(f_2 - \hat{\Pi}_{V_2}f_2)\|_n. \end{aligned} \quad (6.10)$$

From (6.9) and (6.10), we conclude that

$$\|(\hat{\Pi}_V f_2)_1\|_n^2 = \|(\hat{\Pi}_V(f_2 - \hat{\Pi}_{V_2}f_2))_1\|_n^2 \leq \frac{1}{(1 - \rho_{0,\delta}^2)^2} \|\hat{\Pi}_{V_1}(f_2 - \hat{\Pi}_{V_2}f_2)\|_n^2.$$

Combining this with (6.8) and (5.15) gives (6.5). This completes the proof. \square

6.2. End of proof of Theorem 3 and 4. The only place where we modify the proof of Theorem 2 is the analysis of the term

$$\frac{2}{(1-\delta)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|f_1 - (\hat{\Pi}_V f)_1\|_n^2 \right], \quad (6.11)$$

which we bounded by

$$\frac{1+\delta}{(1-\delta)^2} \frac{6}{1-\rho_0^2} (\|f_1 - \Pi_{V_1} f_1\|^2 + \|f_2 - \Pi_{V_2} f_2\|^2), \quad (6.12)$$

by using Lemma 9. We show that, under the additional Assumptions (3.6) or (3.7), one can replace the upper bound (6.12) by the ones given in Theorem 3 and 4, respectively. To achieve this, we replace Lemma 9 by Proposition 4 and 5.

In order to proof Theorem 3 we decompose

$$f_1 - (\hat{\Pi}_V f)_1 = f_1 - (\Pi_V f_1)_1 + (\Pi_V f_1)_1 - (\hat{\Pi}_V f_1)_1 - (\hat{\Pi}_V f_2)_1.$$

Thus

$$\begin{aligned} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|(f_1 - (\hat{\Pi}_V f)_1)\|_n^2 \right] &\leq 3\mathbb{E} [\|f_1 - (\Pi_V f_1)_1\|_n^2] \\ &\quad + 3\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\Pi_V f_1)_1 - (\hat{\Pi}_V f_1)_1\|_n^2 \right] \\ &\quad + 3\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\hat{\Pi}_V f_2)_1\|_n^2 \right]. \end{aligned}$$

The third term on the right-hand side is part of Proposition 5. Using Lemma 1 and the projection theorem, the first term can be bounded by

$$\begin{aligned} \mathbb{E} [\|f_1 - (\Pi_V f_1)_1\|_n^2] &= \|f_1 - (\Pi_V f_1)_1\|^2 \\ &\leq \frac{1}{(1-\rho_0^2)} \|f_1 - \Pi_V f_1\|^2 \\ &\leq \frac{1}{(1-\rho_0^2)} \|f_1 - \Pi_{V_1} f_1\|^2. \end{aligned} \quad (6.13)$$

Applying Proposition 1 and the projection theorem, the second term can be bounded by

$$\begin{aligned}
& \mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\Pi_V f_1)_1 - (\hat{\Pi}_V f_1)_1\|_n^2 \right] \\
& \leq \frac{(1+\delta)}{(1-\delta)} \frac{1}{(1-\rho_0^2)} \mathbb{E} \left[\|\Pi_V f_1 - \hat{\Pi}_V f_1\|_n^2 \right] \\
& \leq \frac{(1+\delta)}{(1-\delta)} \frac{1}{(1-\rho_0^2)} \mathbb{E} \left[\|f_1 - \Pi_V f_1\|_n^2 \right] \\
& \leq \frac{(1+\delta)}{(1-\delta)} \frac{1}{(1-\rho_0^2)} \|f_1 - \Pi_{V_1} f_1\|^2.
\end{aligned}$$

This completes the proof of Theorem 3. In order to proof Theorem 4 we decompose

$$f_1 - (\hat{\Pi}_V f)_1 = f_1 - (\Pi_V f_1)_1 - (\Pi_V f_2)_1 + (\Pi_V f)_1 - (\hat{\Pi}_V f)_1.$$

Thus

$$\begin{aligned}
\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(f_1 - (\hat{\Pi}_V f)_1)\|^2 \right] & \leq 3\mathbb{E} \left[\|f_1 - (\Pi_V f_1)_1\|_n^2 \right] \\
& \quad + 3\mathbb{E} \left[\|(\Pi_V f_2)_1\|_n^2 \right] \\
& \quad + 3\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\Pi_V f)_1 - (\hat{\Pi}_V f)_1\|_n^2 \right].
\end{aligned}$$

The first term on the right-hand side is bounded in (6.13), the second one in Proposition 4. Using the definition of \mathcal{E}_δ and Lemma 1, we obtain

$$\begin{aligned}
\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\Pi_V f)_1 - (\hat{\Pi}_V f)_1\|_n^2 \right] & \leq (1+\delta) \mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\Pi_V f)_1 - (\hat{\Pi}_V f)_1\|^2 \right] \\
& \leq (1+\delta) \frac{1}{(1-\rho_0^2)} \mathbb{E} \left[1_{\mathcal{E}_\delta} \|\Pi_V f - \hat{\Pi}_V f\|^2 \right].
\end{aligned}$$

By Proposition 7, this can be bounded by

$$\frac{(1+\delta)}{(1-\delta)^2} \frac{2}{(1-\rho_0^2)} \frac{\varphi^2 d}{n} (\|f_1 - \Pi_{V_1} f_1\|^2 + \|f_2 - \Pi_{V_2} f_2\|^2).$$

This completes the proof. \square

APPENDIX A. PROOF OF LEMMA 1

We first show how (i) implies (ii) and (iii). Let $h_1 \in \mathcal{H}_1$ and $h_2 \in \mathcal{H}_2$. Then by (i) we have $\|h_1 + h_2\|^2 \geq \|h_1\|^2 - 2\varrho \|h_1\| \|h_2\| + \|h_2\|^2$ and (ii) follows from the inequality $2\varrho \|h_1\| \|h_2\| \leq \|h_1\|^2 + \|h_2\|^2$, while (iii) follows from $2\varrho \|h_1\| \|h_2\| \leq \varrho^2 \|h_1\|^2 + \|h_2\|^2$.

Next, we show how (ii) implies (i). Let $0 \neq h_1 \in \mathcal{H}_1$ and $0 \neq h_2 \in \mathcal{H}_2$. We may assume without loss of generality that $\|h_1\| = \|h_2\| = 1$ and

that $\langle h_1, h_2 \rangle \geq 0$. Then by (ii) we have $2 - 2\langle h_1, h_2 \rangle = \|h_1 - h_2\|^2 \geq 2(1 - \varrho)$ which gives (i).

Finally, suppose that (iii) is true. Let $0 \neq h_1 \in \mathcal{H}_1$ and $0 \neq h_2 \in \mathcal{H}_2$. Again, we may assume that $\|h_1\| = \|h_2\| = 1$. Then by (iii) we have $1 - \langle h_1, h_2 \rangle^2 = \|h_1 - \langle h_1, h_2 \rangle h_2\|^2 \geq 1 - \varrho^2$ which gives (i). This completes the proof. \square

APPENDIX B. A FEASIBLE ESTIMATOR

In this appendix, we show that estimators based on the condition $(1/n) \sum_{i=1}^n g_1(X_1^i) = 0$ have (up to a constant and a term of smaller order) the same risk bound as our estimators based on the condition $\mathbb{E}[g_1(X_1)] = 0$. We only sketch the main arguments in the case $W_1 = V_1$. Suppose that we choose $U_1 \subset L^2(\mathbb{P}^{X_1})$ and $V_2 \subset L^2(\mathbb{P}^{X_2})$, where V_2 contains all constant functions. Let $V'_1 = \{g_1 \in U_1 \mid (1/n) \sum_{i=1}^n g_1(X_1^i) = 0\}$ and $V_1 = \{g_1 \in U_1 \mid \mathbb{E}[g_1(X_1)] = 0\}$. Since V_2 contains all constants, we have $V = V_1 + V_2 = V'_1 + V_2$. This implies that the first components of $\hat{f}_{V_1+V_2}$ and $\hat{f}_{V'_1+V_2}$ in V_1 and V'_1 , respectively, differ only by the constant

$$\frac{1}{n} \sum_{i=1}^n (\hat{f}_V)_1(X_1^i).$$

The risk of this constant can be bounded as follows. By the bound $(x + y)^2 \leq 2x^2 + 2y^2$, we have

$$\begin{aligned} & \mathbb{E} \left[1_{\mathcal{E}_\delta} \left(\frac{1}{n} \sum_{i=1}^n (\hat{f}_V)_1(X_1^i) \right)^2 \right] \\ & \leq 2\mathbb{E} \left[1_{\mathcal{E}_\delta} \left(\frac{1}{n} \sum_{i=1}^n (\hat{f}_V)_1(X_1^i) - f_1(X_1^i) \right)^2 \right] \\ & \quad + 2\mathbb{E} \left[1_{\mathcal{E}_\delta} \left(\frac{1}{n} \sum_{i=1}^n f_1(X_1^i) \right)^2 \right]. \end{aligned}$$

Applying the Cauchy-Schwarz inequality and the fact that the $f_1(X_1^i)$ are independent and centered, this can be bounded by

$$2\mathbb{E} \left[1_{\mathcal{E}_\delta} \|(\hat{f}_V)_1 - f_1\|_n^2 \right] + \frac{2\|f_1\|^2}{n}.$$

Now apply Lemma 9 to the first term.

APPENDIX C. PROOF OF LEMMA 3 AND 4

First, we prove Lemma 3. In Section 4.1, we have shown that

$$\|g_1\|_\infty^2 \leq \varphi_1^2 d_1 \|g_1\|^2$$

and

$$\|g_{j2}\|_\infty^2 \leq \varphi_2^2 d_{j2} \|g_{j2}\|^2$$

for all $g_1 \in V_1$, $g_{j2} \in V_{j2}$, $1 \leq j \leq q-1$, with $\varphi_1^2 = 2/c$ and $\varphi_2^2 = 2/c$. Now let $g = g_1 + g_2 \in V$. Suppose that $g_2 = \sum_{j=1}^{q-1} g_{j2}$ is the decomposition satisfying (4.6). Applying the above bounds, the Cauchy-Schwarz inequality, and Assumption 9, we obtain

$$\|g_2\|_\infty \leq \sum_{j=1}^{q-1} \varphi_2 \sqrt{d_{j2}} \|g_{j2}\| \leq \frac{\varphi_2}{\sqrt{1-\epsilon_2}} \sqrt{\sum_{j=1}^{q-1} d_{j2}} \|g_2\|.$$

Applying again the Cauchy-Schwarz inequality and then Assumption 2 and Lemma 1, we conclude that

$$\begin{aligned} \|g_1 + g_2\|_\infty &\leq \varphi_1 \sqrt{d_1} \|g_1\| + \frac{\varphi_2}{\sqrt{1-\epsilon_2}} \sqrt{\sum_{j=1}^{q-1} d_{j2}} \|g_2\| \\ &\leq \sqrt{\frac{\varphi_1 \vee \varphi_2}{(1-\epsilon_2)(1-\rho_0)}} \sqrt{d_1 + \sum_{j=1}^{q-1} d_{j2}} \|g_1 + g_2\|. \end{aligned} \quad (\text{C.1})$$

This completes the proof of Lemma 3. The proof of Lemma 4 is similar. In [4], it is shown that

$$\|g_1\|_\infty^2 \leq (r_1 + 1)^2 m_1 \int_0^1 g_1^2(x_1) dx_1$$

and

$$\|g_{j2}\|_\infty^2 \leq (r_2 + 1)^2 m_2 \int_0^1 g_{j2}^2(x_{j2}) dx_{j2}$$

for all $g_1 \in V_1$, $g_{j2} \in V_{j2}$, $1 \leq j \leq q-1$. This implies that

$$\|g_1\|_\infty^2 \leq \varphi_1^2 d_1 \|g_1\|^2$$

and

$$\|g_{j2}\|_\infty^2 \leq \varphi_2^2 d_{j2} \|g_{j2}\|^2$$

with $\varphi_1^2 = 2(r_1 + 1)/c$ and $\varphi_2^2 = 2(r_2 + 1)/c$. Now proceed as above. This completes the proof. \square

APPENDIX D. PROOF OF COROLLARY 7

Lemma 10. *Let Assumption 4 and 9 be satisfied. Suppose that (4.17) and (4.18) are satisfied. Then (3.7) is satisfied with*

$$\psi_{\Pi}(V_2) = \sqrt{C_3} \sqrt{\sum_{k=1}^{q-1} d_{2k}^{-2\beta}}$$

and

$$h_1(x_1) = \sqrt{\sum_{k=1}^q h_{1k}^2(x_1)} + \sqrt{\frac{c'}{1 - \epsilon_2} \int \left(\frac{p_X(x_1, x_2)}{p_{X_1}(x_1)p_{X_2}(x_2)} \right)^2 p_{X_2}(x_2) dx_2},$$

where $c' = \sum_{j,k=1, j \neq k}^{q-1} \|h'_{jk}\|^2$. Note that $h_1 \in L^2(\mathbb{P}^{X_1})$, by Assumption 4.

Proof. Let x_1 be fixed (such that $p_X(x_1, \cdot)/(p_{X_1}(x_1)p_{X_2}(\cdot)) \in L^2(\mathbb{P}^{X_2})$, which is satisfied for \mathbb{P}^{X_1} -almost all x_1 , by Assumption 4). By the projection theorem, the expression

$$\int \left(\frac{p_X(x_1, x_2)}{p_{X_1}(x_1)p_{X_2}(x_2)} - g(x_2) \right)^2 p_2(x_2) dx_2,$$

subject to the constraints $g \in H_2$, is minimized by r . Suppose that $r = \sum_{k=1}^{q-1} r_k$ is the decomposition such that (4.6) is satisfied (note that we omit the dependence of r and the r_k on x_1). For $k = 1, \dots, q-1$, we have

$$\mathbb{E} \left[\frac{p_X(x_1, X_2)}{p_{X_1}(x_1)p_{X_2}(X_2)} \middle| X_{2k} = x_{2k} \right] = \frac{p_{X_1, X_{2k}}(x_1, x_{2k})}{p_{X_1}(x_1)p_{X_{2k}}(x_{2k})}. \quad (\text{D.1})$$

Thus the r_k satisfy the $q-1$ equations

$$r_k(x_{2k}) = \frac{p_{X_1, X_{2k}}(x_1, x_{2k})}{p_{X_1}(x_1)p_{X_{2k}}(x_{2k})} - \sum_{j=1, j \neq k}^{q-1} \int r_j(x_{2j}) \frac{p_{X_{2j}, X_{2k}}(x_{2j}, x_{2k})}{p_{X_{2j}}(x_{2j})p_{X_{2k}}(x_{2k})} p_{X_{2j}}(x_{2j}) dx_{2j},$$

for $\mathbb{P}^{X_{2k}}$ -almost all x_{2k} , $1 \leq k \leq q-1$ (note again that we omit the dependence of the r_k on x_1). By (4.17), (4.18), and the Cauchy-Schwarz

inequality, the first and the second term on the right hand side are contained in $\mathcal{H}(\beta, h_{1k}(x_1))$ and $\mathcal{H}(\beta, \sum_{j=1, j \neq k}^{q-1} \|r_j\|_{L^2(\mathbb{P}^{X_{2j}})} \|h'_{jk}\|)$, respectively. We conclude that

$$\begin{aligned} & \|r - \Pi_{V_2} r\|_{L^2(\mathbb{P}^{X_2})} \\ & \leq \sum_{k=1}^{q-1} \|r_k - \Pi_{V_{2k}} r_k\|_{L^2(\mathbb{P}^{X_{2k}})} \\ & \leq \sum_{k=1}^{q-1} C_3 \left(h_{1k}(x_1) + \sum_{j=1, j \neq k}^{q-1} \|r_j\|_{L^2(\mathbb{P}^{X_{2j}})} \|h'_{jk}\| \right) d_{2k}^{-\beta}. \end{aligned}$$

Applying the Cauchy-Schwarz inequality and Assumption 9, this is bounded by

$$\leq \sum_{k=1}^{q-1} C_3 d_{2k}^{-\beta} \left(h_{1k}(x_1) + \sqrt{\frac{\|r\|_{L^2(\mathbb{P}^{X_2})}^2}{1 - \epsilon_2} \sum_{j=1, j \neq k}^{q-1} \|h'_{jk}\|^2} \right).$$

Applying the Cauchy-Schwarz inequality again and the fact that orthogonal projections lower the norm, we obtain the claimed $\psi_{\Pi}(V_2)$ and $h_1(x_1)$. This completes the proof. \square

APPENDIX E. PROOF OF LEMMA 6

We only proof (iii), since (i) and (ii) are standard. By the spectral theorem, there exists an orthogonal matrix V and nonnegative real numbers $\lambda_1(B), \dots, \lambda_{k_1}(B)$ such that

$$B = V^T \text{diag}(\lambda_1(B), \dots, \lambda_{k_1}(B))V. \quad (\text{E.1})$$

Now, by the Cauchy-Schwarz inequality, each entry of a matrix is bounded by the operator norm of that matrix. In particular, we have $|(VAV^T)_{jk}| \leq \|VAV^T\|_{\text{op}} = \|A\|_{\text{op}}$ for all j, k , since V is orthogonal. Applying (E.1), part (i) of this Lemma, the fact that the $\lambda_j(B)$ are nonnegative, and the previous argument, we obtain

$$\begin{aligned} |\text{tr}(AB)| &= \left| \sum_{j=1}^{k_1} (VAV^T)_{jj} \lambda_j(B) \right| \\ &\leq \max_{j=1, \dots, k_1} |(VAV^T)_{jj}| \text{tr}(B) \leq \|A\|_{\text{op}} \text{tr}(B). \end{aligned}$$

This completes the proof. \square

APPENDIX F. PROOF OF (5.16)

In this appendix, we prove (5.16). As mentioned in the proof of Theorem 1, the main arguments are taken from [2, page 139 and 140]. We define the event $\mathcal{A} = \{\|\hat{f}_1\|_\infty \leq k_n\}$. Then

$$\begin{aligned} \mathbb{E} \left[\|f_1 - \hat{f}_1^*\|^2 \right] &= \mathbb{E} \left[(1_{\mathcal{E}_\delta} 1_{\mathcal{A}} + 1_{\mathcal{E}_\delta} 1_{\mathcal{A}^c} + 1_{\mathcal{E}_\delta^c}) \|f_1 - \hat{f}_1^*\|^2 \right] \\ &\leq \mathbb{E} \left[1_{\mathcal{E}_\delta} \|f_1 - \hat{f}_1\|^2 \right] + \mathbb{E} \left[1_{\mathcal{E}_\delta} 1_{\mathcal{A}^c} \|f_1\|^2 \right] + \mathbb{E} \left[1_{\mathcal{E}_\delta^c} (\|f_1\| + k_n)^2 \right]. \end{aligned}$$

Thus it remains to consider the last two expressions. By Theorem 7, the last one is bounded by

$$2^{3/4} (\|f_1\| + k_n)^2 d \exp \left(-\kappa \frac{n\delta^2}{\varphi^2 d} \right).$$

Consider the other one. By Assumption 3, we have $\|\hat{f}_1\|_\infty^2 \leq \varphi^2 d \|\hat{f}_1\|^2$. If \mathcal{E}_δ holds, then

$$\|\hat{f}_1\|_\infty^2 \leq \frac{\varphi^2 d}{(1-\delta)} \|\hat{\Pi}_{W_1}(\hat{\Pi}_V \mathbf{Y})_1\|_n^2 \leq \frac{\varphi^2 d}{(1-\delta)} \|(\hat{\Pi}_V \mathbf{Y})_1\|_n^2,$$

where we applied the definition of \mathcal{E}_δ and the fact that projections lower the norm. By Proposition 1, the last expression is bounded by

$$\frac{(1+\delta)\varphi^2 d}{(1-\delta)^2(1-\rho_0^2)} \|\hat{\Pi}_V \mathbf{Y} - g_2\|_n^2,$$

for $g_2 \in V_2$ arbitrary. Using $\|\hat{\Pi}_V \mathbf{Y} - g_2\|_n \leq \|\hat{\Pi}_V(f - g_2)\|_n + \|\hat{\Pi}_V \epsilon\|_n \leq \|f - g_2\|_n + \|\epsilon\|_n$ and Markov's inequality, we conclude that

$$\begin{aligned} \mathbb{P}(\mathcal{E}_\delta \cap \mathcal{A}^c) &\leq \mathbb{P} \left(\frac{(1+\delta)\varphi^2 d}{(1-\delta)^2(1-\rho_0^2)} (\|f - g_2\|_n + \|\epsilon\|_n)^2 > k_n^2 \right) \\ &\leq \frac{2(1+\delta)\varphi^2 d (\|f - g_2\|^2 + \sigma^2)}{(1-\delta)^2(1-\rho_0^2)k_n^2}. \end{aligned}$$

Letting $g_2 = \Pi_{V_2} f$, this completes the proof. \square

ACKNOWLEDGEMENTS

Finally, I sincerely would like to thank Prof. Enno Mammen for his support and guidance during the preparation of this paper.

REFERENCES

- [1] N. Aronszajn. Theory of reproducing kernels. *Trans. Amer. Math. Soc.*, 68:337–404, 1950.
- [2] Y. Baraud. Model selection for regression on a random design. *ESAIM Probab. Statist.*, 6:127–146, 2002.

- [3] P. J. Bickel, C. A. J. Klaassen, Y. Ritov, and J. A. Wellner. *Efficient and adaptive estimation for semiparametric models*. Springer, New York, 1998. Reprint of the 1993 original.
- [4] L. Birgé and P. Massart. From model selection to adaptive estimation. In *Festschrift for Lucien Le Cam*, pages 55–87. Springer, New York, 1997.
- [5] L. Birgé and P. Massart. Minimum contrast estimators on sieves: exponential bounds and rates of convergence. *Bernoulli*, 4:329–375, 1998.
- [6] L. Breiman and J. H. Friedman. Estimating optimal transformations for multiple regression and correlation. *J. Amer. Statist. Assoc.*, 80:580–598, 1985.
- [7] A. Dalalyan, Y. Ingster, and A. B. Tsybakov. Statistical inference in compound functional models. *Probab. Theory Related Fields*, 158:513–532, 2014.
- [8] C. De Boor. *A practical guide to splines*. Springer, New York, revised edition, 2001.
- [9] R. M. Dudley. *Real analysis and probability*. Cambridge University Press, Cambridge, 2002. Revised reprint of the 1989 original.
- [10] S. Efromovich. Nonparametric regression with the scale depending on auxiliary variable. *Ann. Statist.*, 41:1542–1568, 2013.
- [11] J. Fan, W. Härdle, and E. Mammen. Direct estimation of low-dimensional components in additive models. *Ann. Statist.*, 26:943–971, 1998.
- [12] T. J. Hastie and R. J. Tibshirani. *Generalized additive models*. Chapman and Hall, London, 1990.
- [13] J. Horowitz, J. Klemelä, and E. Mammen. Optimal estimation in additive regression models. *Bernoulli*, 12:271–298, 2006.
- [14] J. Huang, J. L. Horowitz, and F. Wei. Variable selection in nonparametric additive models. *Ann. Statist.*, 38:2282–2313, 2010.
- [15] S. Kayalar and H. L. Weinert. Error bounds for the method of alternating projections. *Math. Control Signals Systems*, 1:43–59, 1988.
- [16] V. Koltchinskii and M. Yuan. Sparsity in multiple kernel learning. *Ann. Statist.*, 38:3660–3695, 2010.
- [17] H. O. Lancaster. Some properties of the bivariate normal distribution considered in the form of a contingency table. *Biometrika*, 44:289–292, 1957.
- [18] P. D. Lax. *Functional analysis*. John Wiley and Sons, Inc. New York, 2002.
- [19] O. B. Linton. Efficient estimation of additive nonparametric regression models. *Biometrika*, 84:469–473, 1997.
- [20] E. Mammen, O. Linton, and J. Nielsen. The existence and asymptotic properties of a backfitting projection algorithm under weak conditions. *Ann. Statist.*, 27:1443–1490, 1999.
- [21] L. Meier, S. van de Geer, and P. Bühlmann. High-dimensional additive modeling. *Ann. Statist.*, 37:3779–3821, 2009.
- [22] A. Muro and S. van de Geer. The additive model with different smoothness for the components. <http://arxiv.org/abs/1405.6584>, 2014.
- [23] J. D. Opsomer. Asymptotic properties of backfitting estimators. *J. Multivariate Anal.*, 73:166–179, 2000.
- [24] J. D. Opsomer and D. Ruppert. Fitting a bivariate additive model by local polynomial regression. *Ann. Statist.*, 25:186–211, 1997.
- [25] G. Raskutti, M. J. Wainwright, and B. Yu. Minimax-optimal rates for sparse additive models over kernel classes via convex programming. *J. Mach. Learn. Res.*, 13:389–427, 2012.

- [26] H. Rauhut. Compressive sensing and structured random matrices. In *Theoretical Foundations and Numerical Methods for Sparse Recovery*, Radon Ser. Comput. Appl. Math., 9, pages 1–92. Walter de Gruyter, Berlin, 2010.
- [27] M. Reed and B. Simon. *Methods of modern mathematical physics. I. Functional analysis*. Academic Press, Inc., New York, 2nd edition, 1980.
- [28] M. Rudelson. Random vectors in the isotropic position. *J. Funct. Anal.*, 164:60–72, 1999.
- [29] M. Rudelson and R. Vershynin. Sampling from large matrices: an approach through geometric functional analysis. *J. ACM*, 54:19 pp, 2007.
- [30] C. J. Stone. Additive regression and other nonparametric models. *Ann. Statist.*, 13:689–705, 1985.
- [31] T. Suzuki and M. Sugiyama. Fast learning rate of multiple kernel learning: trade-off between sparsity and smoothness. *Ann. Statist.*, 41:1381–1405, 2013.
- [32] A. B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, New York, 2009.
- [33] S. A. Van de Geer. *Applications of empirical process theory*. Cambridge University Press, Cambridge, 2000.
- [34] J. Von Neumann. *Functional Operators. II. The Geometry of Orthogonal Spaces*. Princeton University Press, Princeton, 1950.
- [35] J. Weidmann. *Lineare Operatoren in Hilberträumen. Teil 1. Grundlagen*. B. G. Teubner, Stuttgart, 2000.

INSTITUT FÜR ANGEWANDTE MATHEMATIK, UNIVERSITÄT HEIDELBERG, IM
NEUENHEIMER FELD 294, 69120 HEIDELBERG, GERMANY
E-mail address: wahl@uni-heidelberg.de